

LANE: LABEL-AWARE NOISE ELIMINATION FOR FINE-GRAINED TEXT CLASSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we propose Label-Aware Noise Elimination (LANE), a new approach that improves the robustness of deep learning models when trained under increased label noise in fine-grained text classification. LANE leverages the semantic relations between classes and monitors the training dynamics of the model on each training example to dynamically lower the importance of training examples that are perceived to have noisy labels. We test the effectiveness of LANE in fine-grained text classification and benchmark our approach on a wide variety of datasets with various number of classes and various amounts of label noise. LANE considerably outperforms strong baselines on all datasets, obtaining significant improvements ranging from an average improvement of 2.4% in F1 on manually annotated datasets to a considerable average improvement of 4.5% F1 on datasets with higher levels of label noise. We carry out comprehensive analyses of LANE and identify the key components that lead to its success.

1 INTRODUCTION

Deep learning models are increasingly powerful in many NLP applications, but their success is often hindered by data quality. Many existing datasets are annotated by humans on crowdsourcing platforms Demszky et al. (2020) or by automatic approaches such as distant (or weak) supervision Mintz et al. (2009); Wang et al. (2012); Abdul-Mageed & Ungar (2017), and, while weak supervision inherently introduces unwanted mislabeled examples, humans—no matter how careful, are also prone to making labeling errors especially on fine-grained tasks that involve distinguishing between a large number of closely confusable or overlapping classes, e.g., emotion detection Mohammad (2012); Islam et al. (2019); Bao et al. (2009); Strapparava et al. (2012); Liu et al. (2019) or fine-grained topic classification tasks Lewis et al. (2004). The mislabeled training examples are particularly harmful when learning large overparameterized neural networks, since these networks can achieve zero training error on any dataset, with very poor generalization capabilities Zhang et al. (2016).

Several works Li et al. (2023); Karim et al. (2022); Liu & Guo (2020) designed various changes to the training process to learn under label noise. For example, Peer Loss Function Liu & Guo (2020) alters the training loss function to account for label noise, DISC Li et al. (2023) utilizes an instance-specific dynamic thresholding mechanism that blocks access to specific training examples based on the momentum of each instance’s memorization strength. Unicon Karim et al. (2022) leverages a semi-supervised learning (SSL) framework that considers potentially noisy labeled data as unlabeled examples in an SSL algorithm. Area Under the Margin (AUM) Pleiss et al. (2020) utilizes an instance-specific average margin that identifies potentially mislabeled examples from the training set according to the model’s behavior on these examples and blocks access to these examples through a fixed threshold. AUM measures the average difference between the logit values corresponding to a sample’s *assigned* label and its largest *non-assigned* label calculated across the training epochs. The AUM for a *mislabeled sample* is expected to be low, likely negative since the model—through generalization from other correctly labeled training samples, tends to predict the sample in its (hidden) true class which is different from the (wrongly) assigned class, and hence, the largest logit (among all logits) no longer corresponds to the assigned (wrong) label Pleiss et al. (2020). After this data characterization by AUM, Pleiss et al. (2020) subsequently remove samples with low AUM from the training set using a fixed rigid AUM threshold (i.e., the 95 percentile).

054 However, we posit that, through this fixed threshold used to remove mislabeled samples, difficult but
055 valuable samples that exist under the threshold are unnecessarily removed from the training set. In
056 addition, the computation of AUM that contrasts two labels (the assigned—potentially wrong—label
057 and the largest non-assigned label) treats labels independently, and thus, ignores semantic similarities
058 that inherently exist between fine-grained classes (e.g., in fine-grained emotion detection tasks, “anger”
059 is semantically more similar to “fear” than it is to “joy”). To this end, we introduce Label-Aware
060 Noise Elimination (LANE), a novel approach that identifies mislabeled or noisy samples from the
061 training data and seamlessly mitigates their harmful effects. Unlike Pleiss et al. (2020) who remove
062 mislabeled or ambiguous samples from the training set using a fixed threshold, we improve the
063 robustness of our model under label noise by retaining *all* training samples but re-weighting them
064 differently based on the model’s behavior on these samples measured against their assigned labels. In
065 re-weighting the samples, we estimate the degree of “noisiness” of the assigned labels by introducing
066 *label-aware margins* averaged across training iterations that capture inter-class semantic similarities.
067 For example, a sample with true label “anger” but with assigned label “joy” is noisier (has a higher
068 degree of noisiness) than a sample with true label “anger” but with assigned label “fear” since “fear”
069 is semantically closer to “anger” than “joy”. Our label-aware margins extend the concept of *margins*
070 Pleiss et al. (2020) by adaptively weighting samples when the (hidden) true label and the (wrongly)
071 assigned label do not match. Precisely, we capture inter-class semantic similarities and dynamically
072 lower samples’ weights if the model perceives them as noisy (the noisier the assigned label the lower
073 the weight). We learn the inter-class semantic similarities using a label-aware supervised contrastive
074 loss to improve the capabilities of the model to distinguish between easily confusable samples by
075 bringing the latent representations of input samples closer together if they belong to semantically
076 similar classes and further apart if they belong to semantically dissimilar classes.

076 We evaluate the effectiveness of LANE on multiple well-established fine-grained datasets: Empathetic
077 Dialogues Rashkin et al. (2019), GoEmotions Demszky et al. (2020), ISEAR Scherer & Wallbott
078 (1994), CancerEMO Sosea & Caragea (2020), RCV1 Lewis et al. (2004), SciHTC Sadat & Caragea
079 (2022), SST-5 Socher et al. (2013a), Amazon Review McAuley & Leskovec (2013), Yelp Review
080 Asghar (2016), and Yahoo Answer Chang et al. (2008). Using these datasets, we show that LANE
081 works well on various tasks and domains (emotion and general text classification; social networks,
082 dialogues, and personal experiences). In all our experiments, automatically scaling down the impor-
083 tance of identified noisy samples from the training set shows great potential, improving the overall
084 performance on our original datasets by 2.4% F1 on average over the strong AUM approach Pleiss
085 et al. (2020) and by 4.5% F1 on average on our datasets with higher levels of label noise.

086 We summarize our contributions as follows: **1)** We introduce LANE, a new approach that leverages
087 inter-class semantic similarities and monitors the training dynamics of each training example to
088 automatically identify and minimize the harmful effects of ambiguous or mislabeled examples; **2)**
089 We evaluate the effectiveness of our approach on ten text classification benchmark datasets from
090 different tasks and domains; **3)** We carry out a comprehensive analysis and ablation study of LANE
091 and analyze how it performs on datasets that have different levels of noise.

092 2 RELATED WORK

095 Learning with label noise have started to received substantial attention due to the high risk of deep
096 learning models to overfit Liu & Tao (2015); Goldberger & Ben-Reuven (2016); Ren et al. (2018);
097 Englesson & Azizpour (2021); Zhang & Plank (2021); Margatina et al. (2021); Li et al. (2021); Plank
098 (2022); Karim et al. (2022); Garg et al. (2023); Wei et al. (2023c;b;a). For example, Goldberger &
099 Ben-Reuven (2016) propose adding a noise layer in the neural network architecture, whose parameters
100 can be learned for an accurate label estimation. Saxena et al. (2019) introduce a curriculum-learning
101 approach that uses learnable data parameters to rank the importance of examples in the learning
102 process. These parameters are then leveraged to decide the data to use at different training stages.
103 Liu & Guo (2020) on the other hand propose to alter the loss function to make it more robust in
104 the face of label noise and introduce Peer Loss Functions, which evaluate predictions on both the
105 samples at hand, as well as carefully automatically constructed *peer* samples. Other approaches focus
106 on data quality and design techniques to accurately identify and eliminate potentially mislabeled
107 instances Brodley & Friedl (1999); Pleiss et al. (2020); Swayamdipta et al. (2020). For example,
Swayamdipta et al. (2020) introduce data cartography, a model-based tool that separates training data
into three (potentially overlapping) regions, easy-to-learn, ambiguous, and hard-to-learn (many of

108 which are mislabeled), and re-trains on each data region to understand its benefits to learning and
 109 generalization. Pleiss et al. (2020) identify and subsequently remove mislabeled training samples
 110 by monitoring the behavior of the model on each sample and estimating its Area Under the Margin
 111 (AUM) to determine what to remove from the data. Our work builds on this approach: we reformulate
 112 the Area Under the Margin Pleiss et al. (2020) and leverage the inter-class semantic similarities
 113 present in fine-grained tasks to improve training data quality and diminishing the harmful effects of
 114 noisy samples by reweighting the importance of samples during training.

115 The idea of weighting each training example has been well studied in the literature. A classical method
 116 in statistics is importance sampling Kahn & Marshall (1953), which assigns weights to samples in
 117 order to align one distribution to another. Boosting algorithms such as AdaBoost Freund et al. (1999),
 118 select harder samples to train subsequent classifiers. Focal loss Lin et al. (2017) incorporates a soft
 119 weighting scheme that puts emphasis on harder samples. Similarly, hard sample mining Shrivastava
 120 et al. (2011) reduces samples in the majority class and selects the most difficult samples to perform
 121 training on. In contrast to these works, our weighting mechanism exploits the similarities between
 122 classes and ensures noisy samples do not play a significant part in model training.

123 Supervised contrastive learning is an approach that brings the latent representations of input samples
 124 closer together if they belong to the same class (*positives*) and further apart if they belong to different
 125 classes (*negatives*). Gunel et al. (2020) use a supervised contrastive loss to improve fine-tuning
 126 performance of pre-trained language models in several few-shot learning scenarios. Khosla et al.
 127 (2020) introduce a variation of the traditional contrastive loss which aims to produce more samples in
 128 the *positive* set. Instead of only considering samples with the same class as belonging to the positive
 129 set, they propose to use data augmentation to generate more positive samples. Suresh & Ong (2021)
 130 build upon this approach but argue that not all negative samples are equal. To this end, they propose
 131 Label-aware Contrastive Loss (LCL) that learns a weight network to infer the relations between
 132 classes and weigh samples differently. In contrast, our LANE approach exploits *label-aware margins*
 133 to improve the robustness under label noise.

134 3 PROPOSED APPROACH

135 Here, we first provide background on Area Under the Margin introduced by Pleiss et al. (2020) (§3.1)
 136 and then present Label-Aware Noise Elimination (LANE), our new approach that improves model
 137 robustness in the face of label noise (§3.2).

138 3.1 BACKGROUND

139 Area Under the Margin (AUM) Pleiss et al. (2020) is a well-established approach that monitors the
 140 training dynamics of examples by analyzing their margins during training epochs to automatically
 141 identify and remove mislabeled examples from the training data. At training epoch t , the margin M
 142 Pleiss et al. (2020); Bartlett et al. (2017); Elsayed et al. (2018); Jiang et al. (2018) of an example \mathbf{x}
 143 with assigned label y is defined as follows:
 144

$$145 M^{(t)}(\mathbf{x}, y) = z_y^{(t)}(\mathbf{x}) - \max_{i \neq y} z_i^{(t)}(\mathbf{x}) \quad (1)$$

146 where $z_y^{(t)}(\mathbf{x})$ is the logit corresponding to assigned label y , and $\max_{i \neq y} z_i^{(t)}(\mathbf{x})$ is the largest *other*
 147 logit corresponding to label i (from among all non-assigned labels). The margin measures how
 148 different the assigned label is compared to a model’s *belief* in a label at some epoch. A negative
 149 margin likely implies an incorrect prediction, whereas a positive margin implies a correct prediction.
 150 The contribution to generalization of an example \mathbf{x} is measured by averaging the margins of \mathbf{x} across
 151 all training epochs T which represents the Area Under the Margin (AUM) Pleiss et al. (2020), defined
 152 as follows:
 153

$$154 \text{AUM}(\mathbf{x}, y) = \frac{1}{T} \sum_{t=1}^T M^{(t)}(\mathbf{x}, y) \quad (2)$$

155 Figure 1 shows the AUMs of two examples (one correctly labeled and another incorrectly labeled)
 156 from an emotion dataset. In the first example, *Makes me sad how brain damage affects boxers*
 157

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

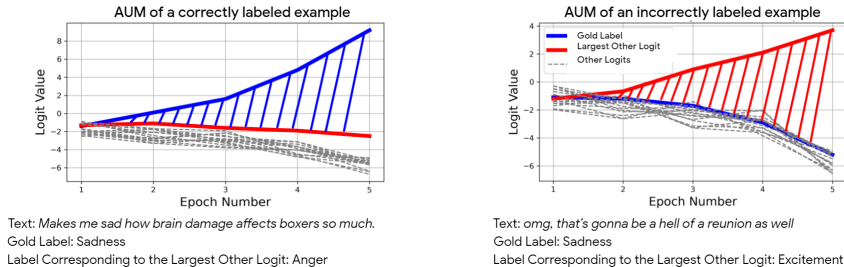


Figure 1: Comparison between the AUM of a correctly labeled example and a mislabeled example.

so much, its assigned (or gold) label is “sadness” which is correct and we observe how the logit corresponding to the assigned label grows larger in each epoch, resulting in a positive high AUM. In contrast, in the second example *omg, that’s gonna be a hell of a reunion as well*, the assigned (gold) label is “sadness”, which is unarguably incorrect, and we observe how other logits, such as the logit corresponding to the “excitement” emotion, are consistently larger than the logit of the “sadness” emotion since the model learns through generalization (from other training examples) that this example shares characteristics of the “excitement” class. Consequently, this example has a low AUM, indicating that its assigned (gold) label is noisy.

Pleiss et al. (2020) first identify mislabeled samples by learning a threshold of separation between the AUMs of clean and erroneous samples through a new artificial class that mimics the training dynamics of mislabeled data and then remove all samples that fall under this threshold. We identified two limitations of the AUM approach. First, we observed (through manual inspection) that through this fixed threshold elimination, difficult but valuable samples that fall under the threshold are unnecessarily removed, and hence, the model has access to less diverse and challenging samples. Second, the current formulation of AUM considers a uniform penalty for each mislabeled sample, irrespective of the semantic similarity between fine-grained classes. A mislabeled example should have a larger negative margin when the wrongly assigned label is more distant from the (hidden) true label and a smaller negative margin when the wrongly assigned label is closer to the (hidden) true label. For example, a sample expressing “excitement” (hidden true label) should have a larger negative margin if the sample is wrongly annotated as “sadness” and a smaller negative margin if the sample is wrongly annotated as “joy”. Thus, we argue that the margin M should take into account the inter-class semantic similarities and incur a higher penalty for semantically distant classes and a lower penalty for closely related classes.

3.2 OUR PROPOSAL: LABEL-AWARE NOISE ELIMINATION

We now introduce LANE, our new approach that addresses the above limitations and improves model robustness on fine-grained text classification under label noise. In our approach we redefine the concept of margin to *label-aware margin* to account for the inter-class semantic similarities. Moreover, instead of unnecessarily removing difficult but valuable samples from the training set if they fall under the fixed AUM threshold, we use all samples from the training data, however weighted according to their label-aware margins to reflect inter-class semantic similarities.

Label-aware Margin (LM) Let θ be a classifier that is trained to predict a task (e.g., sentiment analysis) and Π be a weighting network that learns the semantic similarities between classes. To leverage the inherent semantic similarities between classes for dynamic penalty estimation when the assigned label and the prediction do not match we learn a soft-assignment of input samples into all available classes C that accounts for inter-class semantic similarities. Concretely, Π optimizes the following label-aware supervised contrastive loss (learned jointly with our classifier θ):

$$\mathcal{L}_{LSCL} = \sum_{\mathbf{x} \in B} \frac{-1}{|P_{\mathbf{x}}|} \sum_{p \in P_{\mathbf{x}}} \log \frac{w_{\mathbf{x}, y_{\mathbf{x}}} \cdot \exp(h_{\mathbf{x}} \cdot h_p)}{\sum_{k \in B \setminus \{\mathbf{x}\}} w_{\mathbf{x}, y_k} \cdot \exp(h_{\mathbf{x}} \cdot h_k)} \quad (3)$$

where B is the current batch, $P_{\mathbf{x}}$ is the set of positives p for example \mathbf{x} (i.e., in the context of supervised contrastive learning the positives are all examples that belong to the same class as \mathbf{x} and its augmentation Gunel et al. (2020); Khosla et al. (2020)). $h_{\mathbf{x}}$ is the embedding of \mathbf{x} produced by our model θ . $w_{\mathbf{x}, y_{\mathbf{x}}}$ and $w_{\mathbf{x}, y_k}$ represent the soft-assignment of example \mathbf{x} to its assigned label $y_{\mathbf{x}}$ and all

	TEXT	LOGITS								M	LM
		SDN	JOY	FER	ANG	SRP	DSG	TRS	ANT		
\mathbf{x}_1	The doctors do not have any options for him.	1.1	0.45	1.2	1.8	0.27	1.56	0.11	-0.7	-0.6	-0.67
\mathbf{x}_2	I have found so much info and support on this site, and yet they accept me for who I am.	1.1	1.56	1.2	0.45	0.27	0.11	1.8	-0.7	-0.6	-1.15

Table 1: Comparison of Margin (M) and Label-aware Margin (LM) for two examples. The assigned label (fear) is shown in **red bold** and the model predicted label for each example is shown in **blue bold**. For both examples, we observe that M is -0.6 (i.e., $1.2 - 1.8$). In the first example, LM is rescaled slightly since the assigned emotion fear is semantically close to the emotion corresponding to the largest other logit (i.e., anger). In contrast, we observe that in the second example, the assigned emotion fear is semantically distant from the emotion corresponding to the largest other logit which is trust, and hence, LM becomes much smaller.

the other non-assigned labels y_k where $k = 1, \dots, C; k \neq y_{\mathbf{x}}$. To obtain these soft-assignments we utilize the weighting network Π applied on top of our model, where Π can be viewed as a regular linear layer that projects $h_{\mathbf{x}}$ into a vector $\pi_{\mathbf{x}}$ of length C , $\pi_{\mathbf{x}} = \Pi(h_{\mathbf{x}})$. Concretely, $w_{\mathbf{x},y} = \frac{\exp(\pi_{\mathbf{x},y})}{\sum_{i=1}^C \exp(\pi_{\mathbf{x},i})}$.

Using these weights, we propose to rescale the margin and introduce the Label-aware Margin (LM):

$$\text{LM}^{(t)}(\mathbf{x}, y) = \begin{cases} \frac{1}{w_{\mathbf{x},j}} \cdot \text{M}^{(t)}(\mathbf{x}, y) & \text{if } \text{M}^{(t)}(\mathbf{x}, y) < 0 \text{ and } j = \text{argmax}_{i \neq y} z_i^{(t)}(\mathbf{x}) \\ \text{M}^{(t)}(\mathbf{x}, y) & \text{otherwise} \end{cases} \quad (4)$$

where $w_{\mathbf{x},j}$ is the weight obtained using the weighting network Π , which produces higher values if the (potentially wrong) assigned label y of \mathbf{x} is semantically close to the (hidden) true label j predicted by the model, and lower values otherwise (i.e., if the potentially wrong assigned label is semantically distant from the model prediction). Note that we scale the margins only if the margins are negative, since these are the potentially problematic examples that may be overly ambiguous or mislabeled. To showcase the difference between our proposed label-aware margin LM and the vanilla margin M, we show in Table 1 two examples from an emotion dataset alongside the logits produced by the model as well as the margin M and label-aware margin LM. Both of these examples have the assigned label the *fear* emotion—while \mathbf{x}_1 can be viewed as ambiguous, \mathbf{x}_2 is clearly mislabeled. However, although the margin of both examples is the same $M = -0.6$, we notice that the assigned label fear is semantically close to the label corresponding to the largest other logit (i.e., anger)—the model prediction in the first example, whereas in the second example, it is semantically distant from the label corresponding to the largest other logit (i.e., trust)—the model prediction. We emphasize that our LM captures this semantic difference between labels. Specifically, we observe that the LM of the first example, where the prediction and the assigned label are semantically close, i.e., anger and fear, is larger than the LM of the second example where the prediction and the assigned label are semantically distant, i.e., trust and fear.

Average Label-aware Margin (ALM) At an arbitrary iteration t we measure the contribution of training examples to learning and generalization by averaging the LMs across the training process, from the beginning up until the current iteration t and obtain the Average Label-aware Margin (ALM) as follows: $\text{ALM}^{(t)}(\mathbf{x}, y) = \frac{1}{t} \sum_{r=1}^t \text{M}^{(r)}(\mathbf{x}, y)$.

Mitigating the harmful effect of mislabeled examples To mitigate the harmful effect of mislabeled or noisy examples, we propose to use a weighted cross entropy loss during training and assign higher weights for high-ALM examples and lower weights otherwise. Let $N^t = \{\mathbf{x}_i \mid \text{ALM}^{(t)}(\mathbf{x}_i, y_i) < 0\}$ be the set of examples that have negative ALMs up until training iteration t and $\text{ALM}(N^t)$ be the distribution of their ALMs. At t , we propose to scale down the loss on examples from N^t for those examples whose ALM is below the mean of the ALM distribution $\text{ALM}(N^t)$. Specifically, we propose to dynamically fit a truncated Gaussian distribution of mean μ_t and variance σ_t at training iteration t . We assign a weight for each example \mathbf{x}_i at iteration t as follows:

$$\lambda_{CE}^t(\mathbf{x}_i, y_i) = \begin{cases} \exp\left(-\frac{(\text{ALM}^{(t)}(\mathbf{x}_i, y_i) - \mu_t)^2}{2\sigma_t^2}\right) & \text{if } \mathbf{x}_i \in N^t \text{ and } \text{ALM}^{(t)}(\mathbf{x}_i, y_i) < \mu_t \\ 1 & \text{otherwise} \end{cases} \quad (5)$$

During training, we estimate the mean μ_t and variance σ_t using the historical predictions of the model:

$$\mu_t = \frac{1}{|N^t|} \sum_{(\mathbf{x}, y) \in N^t} \text{ALM}^{(t)}(\mathbf{x}, y) \quad (6)$$

$$\sigma_t = \frac{1}{|N^t|} \sum_{(\mathbf{x}, y) \in N^t} (\text{ALM}^{(t)}(\mathbf{x}, y) - \mu_t)^2 \quad (7)$$

Intuitively, a low weight for an example indicates that the example produced an ALM that is consistently below the mean of the negative ALM distribution. As we have shown, such examples are potentially mislabeled and may hurt generalization. To mitigate this effect, at each training iteration t we simply rescale the cross entropy loss, assigning lower weight to potentially mislabeled examples:

$$\mathcal{L}_{CE} = \sum_{i=1}^{|B|} \lambda_{CE}^t(\mathbf{x}_i, y_i) \cdot H(\theta(\mathbf{x}_i), y_i) \quad (8)$$

where $\theta(\mathbf{x})$ is the probability distribution of the model θ on example \mathbf{x} , $|B|$ is the batch size, and H is the cross-entropy.

The final loss in LANE is a combination of the weighted cross entropy loss and the contrastive loss:

$$\mathcal{L} = \alpha \cdot \mathcal{L}_{CE} + (1 - \alpha) \cdot \mathcal{L}_{LSCL} \quad (9)$$

In our experiments we set $\alpha = 0.5$.

4 EXPERIMENTS

4.1 LABEL NOISE

We evaluate the effectiveness of LANE on ten datasets under various amounts of label noise. We employ three setups: **1)** Original datasets, where the label noise comes from annotation errors in the dataset collection process, **2)** 20% noise, where we randomly shuffle the labels of 20% of the training data, and **3)** 40% noise, where we perform the same process for 40% of the training examples.

4.2 EXPERIMENTAL SETUP

We carry out all our experiments using an Nvidia A5000 GPU. We use the HuggingFace Transformers Wolf et al. (2020) library for our BERT implementation. The datasets we consider make their train/validation/test splits available, hence, we use the provided splits in our experiments. Similar to Khosla et al. (2020), to expand the positive set of examples in the contrastive loss, we augment our data using synonym replacement Kolomyets et al. (2011), SwitchOut Wang et al. (2018), and backtranslation Tiedemann & Thottingal (2020). In backtranslation we translate from English to German and back to English. For each batch, we generate 7 augmentations. For all datasets we follow the evaluation metrics used in the works introducing the datasets. The initial batch size is set to 32, hence the total batch size (i.e., including augmentations) is 256. In our training setup, we only scale down the importance of examples during training if their ALM is below a threshold that we set as the ALM mean of examples with negative ALMs (Eq. 6). We also experimented with different ALM thresholds such as 0, but observed slightly worse performance than using the mean.

4.3 DATASETS

The datasets used to evaluate LANE are: **1. Empathetic Dialogues** Rashkin et al. (2019), a dataset composed of conversations between a speaker and a listener annotated with 32 emotions. We consider solely the first turn of the conversation in our experiments, resulting in 22,000 total examples. **2. GoEmotions** Demszky et al. (2020), a sentence-level dataset created using Reddit comments that contains more than 58,000 sentences annotated with 27 emotions. **3. ISEAR** (International Survey on Emotion Antecedents and Reactions) Scherer & Wallbott (1994), a dataset of 7,700 personal

324 experiences annotated with 7 emotions. **4. CancerEMO** Sosea & Caragea (2020), a dataset of
 325 8,500 examples collected from a cancer forum annotated at sentence level with the 8 basic Plutchik-8
 326 Plutchik (1980) emotions. **5. RCV1** Lewis et al. (2004), a large scale dataset composed of news
 327 stories labeled with a total of 105 different topics. **6. SciHTC** Sadat & Caragea (2022), a dataset
 328 from 186,160 scientific papers, annotated with 80 possible topics, **7. SST5** Socher et al. (2013b),
 329 a dataset composed of 11,855 sentences from movie reviews, annotated with five sentiment labels:
 330 *negative*, *somewhat negative*, *neutral*, *somewhat positive*, and *positive*. **8. Amazon Review** McAuley
 331 & Leskovec (2013), a sentiment classification dataset composed of 600,000 training and 130,000
 332 test Amazon reviews annotated with 5 sentiment classes. **9. Yelp Review** Asghar (2016), a sentiment
 333 classification dataset with 130,000 training and 10,000 test samples annotated with the same 5
 334 classes, and **10. Yahoo Answer** Chang et al. (2008), a topic classification dataset with 10 topic
 335 classes, composed of 140,000 training and 6,000 test samples.

336 4.4 BASELINE MODELS

337
 338 We use BERT Devlin et al. (2019) base uncased model in all experiments (denoted by BASE). We
 339 compare LANE against methods that use training dynamics to assess the data quality, as well as
 340 approaches focused on exploiting the relationships between classes and approaches aimed at learning
 341 under label noise:

342 **Data Cartography** Following Swayamdipta et al. (2020), we identify three types of training
 343 examples: easy-to-learn (E2L), hard-to-learn (H2L), and ambiguous (AMG) and analyze the
 344 importance of each type to the training process by removing the other two types.

345 **Noise Layer** Following Goldberger & Ben-Reuven (2016), we introduce a noise layer to the BERT
 346 model which we train for correct label estimation. We denote this model by NSE in our experiments.
 347

348 **Peer Loss Function** We also compare our method against Peer Loss Function (PLF) Liu & Guo
 349 (2020), a method that alters the training loss function to account for label noise.

350 **Area Under the Margin** We consider the AUM method Pleiss et al. (2020) as one of our baselines.
 351 This method computes Area Under the Margin metric for each training example and eliminates
 352 low-AUM examples that are potentially noisy, using a fixed threshold for elimination.

353 **Contrastive Learning:** We compare LANE to the label-aware supervised contrastive learning (LCL)
 354 method proposed by Suresh & Ong (2021) and the traditional supervised contrastive learning (SCL)
 355 Khosla et al. (2020).
 356

357 **DISC** Li et al. (2023) proposes an instance-specific dynamic thresholding mechanism that blocks
 358 access to specific training examples based on the momentum of each instance’s memorization
 359 strength. Additionally, DISC proposes to correct the labels of potentially noisy examples.

360 **UNICON** Karim et al. (2022) leverages semi-supervised learning (SSL) to mitigate the harmful
 361 effects of noisy labels by considering the potentially noisy labeled data as unlabeled examples in
 362 an SSL algorithm. UNICON also proposes a new selection mechanism for these unlabeled examples
 363 during training.
 364

365 5 RESULTS

366
 367 **Results on Original Datasets** We show the results on our datasets in Table 2. We make the
 368 following observations. **LANE outperforms the baselines in all setups.** We observe improvements
 369 of 1.6% weighted F1 on ISEAR, 1.4% weighted F1 on RCV1, 1.5% accuracy on Amazon Review
 370 and 1.3% accuracy on Yahoo over the best performing baseline. Notably, over the base BERT model,
 371 we see a 2.9% weighted F1 improvement on GoEmotions and 3% improvement on Yahoo. We note
 372 that LCL, which leverages inter-class relations through the label-aware contrastive learning loss is the
 373 best performing baseline in 5 out of the 10 datasets. Since LANE utilizes similar inter-class relations
 374 during training, we postulate improvements over LCL arise from correctly identifying mislabeled or
 375 ambiguous examples and eliminating their harmful effect during training.

376 **Results on 20% Noise Datasets** The results obtained on the 20% noise (20N) datasets where 20% of
 377 the labels are intentionally flipped are shown in Table 3. We observe that this setup is significantly
 more challenging for the model. For instance, on Empathetic Dialogues the weighted F1 of the BASE

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393

Dataset	Empathetic Dialogues (wF1)	GoEmotions (wF1)	ISEAR (wF1)	CancerEmo (wF1)	RCV1 (wF1)
BASE	58.5 ± 1.2	63.6 ± 1.2	71.5 ± 0.6	75.8 ± 0.8	56.8 ± 0.8
E2L	57.6 ± 0.8	63.2 ± 1.2	71.3 ± 0.7	75.9 ± 0.9	54.3 ± 1.1
H2L	58.9 ± 1.4	64.2 ± 0.7	72.0 ± 0.6	76.3 ± 1.3	55.8 ± 1.4
AMG	59.0 ± 0.6	<u>64.8 ± 0.6</u>	<u>73.4 ± 0.5</u>	76.1 ± 0.8	52.3 ± 1.1
NSE	58.1 ± 1.9	63.8 ± 1.1	72.2 ± 0.8	76.2 ± 0.7	55.7 ± 1.3
PLF	58.4 ± 1.1	63.4 ± 0.8	71.9 ± 1.2	75.9 ± 0.6	56.7 ± 2.2
AUM	58.4 ± 0.6	63.1 ± 1.3	71.8 ± 0.8	76.0 ± 0.9	56.3 ± 0.6
LCL	59.1 ± 1.0	<u>64.8 ± 0.7</u>	72.4 ± 0.5	76.5 ± 0.9	<u>57.9 ± 0.6</u>
SCL	58.9 ± 0.7	62.8 ± 1.1	71.5 ± 0.9	76.2 ± 0.6	56.9 ± 1.7
DISC	<u>59.4 ± 0.9</u>	63.2 ± 1.4	72.3 ± 1.3	76.4 ± 1.1	56.5 ± 1.4
UNICON	58.4 ± 0.7	63.1 ± 0.9	72.5 ± 1.1	<u>76.6 ± 1.3</u>	56.9 ± 1.1
LANE	60.8 ± 0.9	66.5 ± 0.5	74.3 ± 0.4	78.2 ± 0.7	59.3 ± 0.9

DATASET	SciHTC (MF1)	SST-5 (Acc)	Amazon Review (Acc)	Yelp (Acc)	Yahoo (Acc)
BASE	32.5 ± 1.75	56.3 ± 0.6	67.5 ± 0.6	65.9 ± 0.6	75.4 ± 0.6
E2L	31.6 ± 1.5	55.7 ± 1.1	62.9 ± 0.9	62.8 ± 2.3	70.4 ± 1.5
H2L	32.2 ± 1.1	56.6 ± 1.4	67.9 ± 0.8	62.3 ± 1.7	74.1 ± 1.8
AMG	30.6 ± 1.1	55.1 ± 1.3	67.4 ± 1.1	65.1 ± 1.5	72.3 ± 1.7
NSE	32.8 ± 1.5	54.1 ± 1.1	65.8 ± 1.7	65.1 ± 1.3	74.6 ± 1.1
PLF	32.2 ± 1.4	55.7 ± 1.1	67.4 ± 2.1	65.8 ± 1.8	74.8 ± 1.6
AUM	31.2 ± 2.63	56.4 ± 0.9	66.4 ± 0.6	<u>68.1 ± 0.6</u>	72.9 ± 0.6
LCL	<u>33.1 ± 1.42</u>	<u>57.6 ± 0.9</u>	<u>68.2 ± 0.6</u>	66.8 ± 0.6	76.8 ± 0.6
SCL	32.7 ± 1.1	56.8 ± 1.5	67.8 ± 1.3	66.1 ± 1.7	75.3 ± 1.1
DISC	32.8 ± 1.5	56.7 ± 1.3	67.8 ± 2.4	66.4 ± 2.2	75.1 ± 1.7
UNICON	32.7 ± 1.1	56.5 ± 1.6	67.5 ± 1.4	67.9 ± 1.3	<u>77.1 ± 1.5</u>
LANE	34.1 ± 0.87	58.9 ± 0.4	69.7 ± 0.6	69.2 ± 0.6	78.4 ± 0.6

Table 2: Results of LANE on the fine-grained text classification datasets. The reported results are averaged across five runs and standard deviations are provided. Best results are shown in **bold blue** and second best are underlined.

394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411

Dataset	Empathetic Dialogues (wF1)	GoEmotions (wF1)	ISEAR (wF1)	CancerEmo (wF1)	RCV1 (wF1)
BASE	11.6 ± 3.4	21.5 ± 2.8	37.6 ± 3.0	46.7 ± 1.9	44.4 ± 3.8
E2L	10.3 ± 0.8	22.6 ± 1.2	37.1 ± 0.7	47.5 ± 0.9	44.3 ± 1.5
H2L	10.6 ± 1.4	21.8 ± 0.7	37.3 ± 0.6	47.9 ± 1.3	45.8 ± 2.4
AMG	11.4 ± 1.2	22.1 ± 0.6	36.9 ± 0.5	48.4 ± 0.8	45.9 ± 2.7
NSE	10.2 ± 1.9	15.6 ± 1.1	36.4 ± 0.8	44.2 ± 0.7	44.9 ± 1.8
AUM	<u>14.5 ± 0.6</u>	<u>23.5 ± 1.3</u>	38.6 ± 0.8	49.8 ± 0.9	47.6 ± 2.7
SCL	10.4 ± 1.4	21.4 ± 1.3	37.3 ± 0.9	46.4 ± 1.1	45.2 ± 1.5
LCL	10.8 ± 3.24	22.1 ± 5.1	38.3 ± 1.5	46.6 ± 1.2	47.2 ± 2.2
DISC	11.3 ± 1.0	22.5 ± 0.7	40.5 ± 0.5	<u>50.3 ± 0.9</u>	47.1 ± 2.2
UNICON	10.4 ± 1.4	21.9 ± 1.2	39.5 ± 0.9	42.3 ± 0.9	<u>49.2 ± 2.3</u>
LANE	15.9 ± 1.3	24.3 ± 1.2	<u>40.4 ± 0.8</u>	52.5 ± 0.9	49.4 ± 2.1

DATASET	SciHTC (MF1)	SST-5 (Acc)	Amazon Review (Acc)	Yelp (Acc)	Yahoo (Acc)
BASE	24.5 ± 4.6	48.9 ± 3.7	61.5 ± 1.5	60.7 ± 1.3	64.8 ± 1.7
E2L	24.1 ± 2.4	48.2 ± 2.7	60.7 ± 2.4	62.3 ± 2.9	64.9 ± 3.1
H2L	26.7 ± 2.3	48.7 ± 1.9	60.9 ± 2.3	62.6 ± 2.1	65.7 ± 1.8
AMG	26.9 ± 1.4	49.4 ± 1.5	61.3 ± 2.4	62.9 ± 2.3	66.5 ± 1.8
NSE	26.7 ± 4.3	50.4 ± 4.1	61.7 ± 3.5	63.5 ± 3.3	67.2 ± 2.5
AUM	27.4 ± 4.2	50.4 ± 2.5	<u>62.4 ± 1.7</u>	63.3 ± 1.4	65.9 ± 2.4
LCL	24.2 ± 3.9	48.5 ± 5.7	61.7 ± 2.4	63.1 ± 3.1	65.9 ± 3.0
SCL	24.1 ± 3.4	51.5 ± 3.2	62.3 ± 3.5	<u>63.7 ± 3.9</u>	66.8 ± 2.5
DISC	27.5 ± 2.1	<u>51.7 ± 2.6</u>	62.1 ± 2.7	63.2 ± 2.5	<u>67.3 ± 2.1</u>
UNICON	28.9 ± 3.4	50.8 ± 3.1	61.5 ± 3.7	62.3 ± 3.9	64.2 ± 3.7
LANE	30.5 ± 2.97	53.1 ± 1.6	63.1 ± 2.3	65.2 ± 3.1	68.9 ± 2.5

Table 3: Performance of LANE on the ten fine-grained classification datasets in 20% noise setting. The reported results are averaged across five runs and standard deviations are provided. Best results are shown in **bold blue** and second best are underlined.

model drops from 58.5% on the original dataset to 11.6% on the 20N dataset, with a similar trend on all the other datasets. However, even in this more challenging setup, LANE still outperforms all baselines in all setups. For example, on SST5, LANE outperforms AUM in accuracy by 2.7%, DISC by 1.4%, UNICON by 2.3%, and SCL by 1.6%. The improvements over the base model are larger, with an average performance increase of 4.5%.

Results on 40% Noise Datasets We show the results in this high-noise setup in Appendix A.

6 ANALYSIS

425
426
427
428
429
430
431

Ablation Study Here, we analyze the effectiveness of various components of our method. To this end, we first design a version of LANE that uses averaged margins instead of ALMs so that the semantic relations are not incorporated into the model. We achieve this by replacing the ALM term in Eq. 5 with AUM and denote this method by LANE^{-sim}. Second, we investigate the performance of our approach when completely removing the ALM-based weighting. Specifically, we remove the λ weight in Eq. 8 (or set it to 1 always) and train our model to optimize the combination of the contrastive loss and the traditional cross-entropy. We denote this second approach by LANE^{-alm}. Finally, we compare LANE against the vanilla AUM Pleiss et al. (2020), which completely removes

DATASET:	Empathetic Dialogues (wF1)	GoEmotions (mF1)	ISEAR (Acc)	CancerEmo (mF1)	RCV1 (mF1)
LANE ^{-sim}	14.7 ± 1.1	22.9 ± 0.4	39.6 ± 0.5	50.1 ± 0.8	45.2 ± 0.8
LANE ^{-alm}	13.8 ± 0.9	21.6 ± 0.5	37.2 ± 0.8	46.1 ± 0.8	46.2 ± 1.4
AUM	14.5 ± 0.6	23.5 ± 1.3	38.6 ± 0.8	49.8 ± 0.9	47.6 ± 2.7
LANE	15.9 ± 1.3	24.3 ± 1.2	40.4 ± 0.8	52.5 ± 0.9	49.4 ± 2.1
DATASET	SciHTC (MF1)	SST-5 (Acc)	Amazon Review (Acc)	Yelp (Acc)	Yahoo (Acc)
LANE ^{-sim}	28.5 ± 0.8	50.6 ± 0.8	61.2 ± 0.8	62.4 ± 0.8	67.1 ± 0.8
LANE ^{-alm}	29.3 ± <u>1.2</u>	50.2 ± 1.2	61.3 ± 1.2	64.2 ± <u>1.2</u>	66.3 ± 1.2
AUM	27.4 ± 4.2	50.4 ± 2.5	62.4 ± 1.7	63.3 ± 1.4	65.9 ± 2.4
LANE	30.5 ± 2.97	53.1 ± 1.6	63.1 ± 2.3	65.2 ± 3.1	68.9 ± 2.5

Table 4: Ablation study: comparison between LANE, LANE^{-sim}, LANE^{-alm} and vanilla AUM on the datasets using 20% noise. Best results are shown in **bold blue** and second best are underlined.

DATASET:	Empathetic Dialogues (wF1)	GoEmotions (mF1)	ISEAR (ACC)	CancerEMO (mF1)	RCV1 (mF1)
CHATGPT	<u>12.8 ± 3.1</u>	21.4 ± 2.5	37.3 ± <u>1.1</u>	48.9 ± 1.9	42.9 ± 4.6
LLAMA-2	10.9 ± 3.7	20.4 ± 2.7	35.4 ± 1.6	50.2 ± 1.7	39.7 ± 1.8
LANE	15.9 ± 1.3	24.3 ± 1.2	40.4 ± 0.8	52.5 ± 0.9	49.4 ± 2.1
DATASET:	SciHTC (MF1)	SST-5 (Acc)	Amazon Review (Acc)	Yelp (Acc)	Yahoo (Acc)
CHATGPT	28.3 ± 5.0	49.6 ± 0.6	62.6 ± <u>0.9</u>	64.5 ± <u>0.9</u>	64.9 ± 0.9
LLAMA-2	15.1 ± 5.2	54.2 ± 0.4	61.3 ± 2.3	62.3 ± 1.4	61.1 ± 2.3
LANE	30.5 ± 2.97	53.1 ± 1.6	63.1 ± 2.3	65.2 ± 3.1	68.9 ± 2.5

Table 5: Performance of LANE on the ten benchmark datasets compared with LLMs. Best results are shown in **bold blue** and second best are underlined.

examples in the training set that have low AUMs. We show the results obtained on 20N datasets in Table 4. We observe that LANE outperforms LANE^{-sim}, LANE^{-alm} and AUM in all setups. Notably, we see a large improvement on SST-5, where LANE pushes the accuracy score by 2.5% over LANE^{-sim}, by 2.9% over LANE^{-alm} and by 2.6% over AUM. On RCV1, which has a large number of classes, LANE improves the micro F1 score significantly, obtaining 49.4%, a boost of 4.2% over LANE^{-sim}, 3.2% over LANE^{-alm} and 1.8% over AUM. These results show that our proposed Average Label-aware Margin and semantics-aware contrastive loss play an important role in the success of LANE. To gain further insights into LANE we show in Appendix B an error analysis of LANE predictions on the 20% noise ISEAR dataset.

Comparison with Large Language Models We test our approach against few-shot large language models: ChatGPT and Llama-2 13B Touvron et al. (2023) to compare the robustness to label noise of LANE with that of popular LLMs in 20% noise setup. For all datasets except SciHTC we fit a large number of examples in the prompt and set the number of few-shot examples to 100. We use only 10 few-shot examples for SciHTC since the examples (i.e., paper abstracts) are much longer and exceed the context window. Similar to the original 20% noise setup, 20% of the few-shot examples are purposefully mislabeled. To account for the variance produced by the particular few-shot examples selected, we run ChatGPT 10 times with different few-shot examples in the prompt and report average values. Similarly, we run Llama-2 20 times with different few-shot examples and show results in Table 5. We observe that LANE outperforms the LLMs on all datasets except SST5. Notably, LANE improves upon Llama-2 by 15.4% on SciHTC and by 9.7% on RCV1 and improves the performance over ChatGPT by 3.1% accuracy on ISEAR and 6.5% micro F1 on RCV1. Among the LLMs, ChatGPT obtains the best results, outperforming Llama-2 especially in complex tasks such as RCV1 and SciHTC. Concretely, ChatGPT obtains 28.3% macro F1 on RCV1, a 13.2% improvement over Llama-2.

7 CONCLUSION

In this work, we introduced LANE, a new approach that boosts the capabilities of deep learning models when learning under increased label noise. LANE leverages the inter-class semantic similarities and utilizes training dynamics to boost the performance in fine-grained text classification. We tested LANE on ten fine-grained text classification datasets where it obtained improvements in performance over strong baselines and prior works. In the future, we plan to extend our approach to other domains and data types, e.g., image classification and the legal domain. We make our code available to further research in this area.

REFERENCES

- 486
487
488 Muhammad Abdul-Mageed and Lyle Ungar. Emonet: Fine-grained emotion detection with gated
489 recurrent neural networks. In *Proceedings of the 55th Annual Meeting of the Association for*
490 *Computational Linguistics (Volume 1: Long Papers)*, pp. 718–728, 2017.
- 491 Nabiha Asghar. Yelp dataset challenge: Review rating prediction. *arXiv preprint arXiv:1605.05362*,
492 2016.
- 493 S. Bao, S. Xu, L. Zhang, R. Yan, Z. Su, D. Han, and Y. Yu. Joint emotion-topic modeling for social
494 affective text mining. In *2009 Ninth IEEE International Conference on Data Mining*, pp. 699–704,
495 2009.
- 496 Peter L. Bartlett, Dylan J. Foster, and Matus Telgarsky. Spectrally-normalized margin bounds
497 for neural networks. *CoRR*, abs/1706.08498, 2017. URL [http://arxiv.org/abs/1706.](http://arxiv.org/abs/1706.08498)
498 [08498](http://arxiv.org/abs/1706.08498).
- 500 Carla E Brodley and Mark A Friedl. Identifying mislabeled training data. *Journal of artificial*
501 *intelligence research*, 11:131–167, 1999.
- 502 Ming-Wei Chang, Lev-Arie Ratinov, Dan Roth, and Vivek Srikumar. Importance of semantic
503 representation: Dataless classification. In *Aaai*, volume 2, pp. 830–835, 2008.
- 505 Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and
506 Sujith Ravi. GoEmotions: A dataset of fine-grained emotions. In Dan Jurafsky, Joyce Chai, Natalie
507 Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for*
508 *Computational Linguistics*, pp. 4040–4054, Online, July 2020. Association for Computational
509 Linguistics. doi: 10.18653/v1/2020.acl-main.372. URL [https://aclanthology.org/](https://aclanthology.org/2020.acl-main.372)
510 [2020.acl-main.372](https://aclanthology.org/2020.acl-main.372).
- 511 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
512 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*
513 *the North American Chapter of the Association for Computational Linguistics: Human Language*
514 *Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June
515 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL [https:](https://www.aclweb.org/anthology/N19-1423)
516 [//www.aclweb.org/anthology/N19-1423](https://www.aclweb.org/anthology/N19-1423).
- 517 Gamaleldin Fathy Elsayed, Dilip Krishnan, Hossein Mobahi, Kevin Regan, and Samy Bengio. Large
518 margin deep networks for classification. 2018. URL [https://arxiv.org/pdf/1803.](https://arxiv.org/pdf/1803.05598.pdf)
519 [05598.pdf](https://arxiv.org/pdf/1803.05598.pdf).
- 520 Erik Englesson and Hossein Azizpour. Generalized jensen-shannon divergence loss for learning with
521 noisy labels. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan
522 (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 30284–30297. Cur-
523 ran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2021/file/fe2d010308a6b3799a3d9c728ee74244-Paper.pdf)
524 [paper/2021/file/fe2d010308a6b3799a3d9c728ee74244-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/fe2d010308a6b3799a3d9c728ee74244-Paper.pdf).
- 525 Yoav Freund, Robert Schapire, and Naoki Abe. A short introduction to boosting. *Journal-Japanese*
526 *Society For Artificial Intelligence*, 14(771-780):1612, 1999.
- 527 Arpit Garg, Cuong Nguyen, Rafael Felix, Thanh-Toan Do, and Gustavo Carneiro. Instance-dependent
528 noisy label learning via graphical modelling. In *Proceedings of the IEEE/CVF Winter Conference*
529 *on Applications of Computer Vision (WACV)*, pp. 2288–2298, January 2023.
- 530 Jacob Goldberger and Ehud Ben-Reuven. Training deep neural-networks using a noise adaptation
531 layer. 2016.
- 532 Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. Supervised contrastive learning for
533 pre-trained language model fine-tuning. *arXiv preprint arXiv:2011.01403*, 2020.
- 534 Jumayel Islam, Robert E Mercer, and Lu Xiao. Multi-channel convolutional neural network for twitter
535 emotion and sentiment recognition. In *Proceedings of the 2019 Conference of the North American*
536 *Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume*
537 *1 (Long and Short Papers)*, pp. 1355–1365, 2019.
- 538
539

- 540 Yiding Jiang, Dilip Krishnan, Hossein Mobahi, and Samy Bengio. Predicting the generalization gap
541 in deep networks with margin distributions, 2018. URL [https://arxiv.org/abs/1810.](https://arxiv.org/abs/1810.00113)
542 00113.
- 543
- 544 Herman Kahn and Andy W Marshall. Methods of reducing sample size in monte carlo computations.
545 *Journal of the Operations Research Society of America*, 1(5):263–278, 1953.
- 546
- 547 Nazmul Karim, Mamshad Nayeem Rizve, Nazanin Rahnavard, Ajmal Mian, and Mubarak Shah.
548 Unicon: Combating label noise through uniform selection and contrastive learning. In *Proceedings*
549 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9676–9686,
550 June 2022.
- 551 Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron
552 Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural*
553 *information processing systems*, 33:18661–18673, 2020.
- 554
- 555 Oleksandr Kolomiyets, Steven Bethard, and Marie-Francine Moens. Model-portability experiments
556 for textual temporal analysis. In *Proceedings of the 49th Annual Meeting of the Association for*
557 *Computational Linguistics: Human Language Technologies: Short Papers - Volume 2, HLT ’11*,
558 pp. 271–276, USA, 2011. Association for Computational Linguistics. ISBN 9781932432886.
- 559 David D Lewis, Yiming Yang, Tony Russell-Rose, and Fan Li. Rcv1: A new benchmark collection
560 for text categorization research. *Journal of machine learning research*, 5(Apr):361–397, 2004.
- 561
- 562 Xuefeng Li, Tongliang Liu, Bo Han, Gang Niu, and Masashi Sugiyama. Provably end-to-end label-
563 noise learning without anchor points. In Marina Meila and Tong Zhang (eds.), *Proceedings of*
564 *the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine*
565 *Learning Research*, pp. 6403–6413. PMLR, 18–24 Jul 2021. URL [https://proceedings.](https://proceedings.mlr.press/v139/li211.html)
566 [mlr.press/v139/li211.html](https://proceedings.mlr.press/v139/li211.html).
- 567
- 568 Yifan Li, Hu Han, Shiguang Shan, and Xilin Chen. Disc: Learning from noisy labels via dynamic
569 instance-specific selection and correction. In *Proceedings of the IEEE/CVF Conference on*
570 *Computer Vision and Pattern Recognition*, pp. 24070–24079, 2023.
- 571
- 572 Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object
573 detection. In *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988,
574 2017.
- 575
- 576 Chen Liu, Muhammad Osama, and Anderson De Andrade. Dens: A dataset for multi-class emotion
577 analysis. *arXiv preprint arXiv:1910.11769*, 2019.
- 578
- 579 Tongliang Liu and Dacheng Tao. Classification with noisy labels by importance reweighting. *IEEE*
580 *Transactions on pattern analysis and machine intelligence*, 38(3):447–461, 2015.
- 581
- 582 Yang Liu and Hongyi Guo. Peer loss functions: Learning from noisy labels without knowing noise
583 rates. In *International Conference on Machine Learning*, pp. 6226–6236. PMLR, 2020.
- 584
- 585 Katerina Margatina, Giorgos Vernikos, Loïc Barrault, and Nikolaos Aletras. Active learning by
586 acquiring contrastive examples. *arXiv preprint arXiv:2109.03764*, 2021.
- 587
- 588 Julian McAuley and Jure Leskovec. Hidden factors and hidden topics: understanding rating dimen-
589 sions with review text. In *Proceedings of the 7th ACM conference on Recommender systems*, pp.
590 165–172, 2013.
- 591
- 592 Mike Mintz, Steven Bills, Rion Snow, and Daniel Jurafsky. Distant supervision for relation extraction
593 without labeled data. In Keh-Yih Su, Jian Su, Janyce Wiebe, and Haizhou Li (eds.), *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, pp. 1003–1011, Suntec, Singapore,
August 2009. Association for Computational Linguistics. URL [https://aclanthology.](https://aclanthology.org/P09-1113)
[org/P09-1113](https://aclanthology.org/P09-1113).

- 594 Saif Mohammad. #emotional tweets. In **SEM 2012: The First Joint Conference on Lexical and*
595 *Computational Semantics – Volume 1: Proceedings of the main conference and the shared task,*
596 *and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval*
597 *2012)*, pp. 246–255, Montréal, Canada, 7–8 June 2012. Association for Computational Linguistics.
598 URL <https://www.aclweb.org/anthology/S12-1033>.
- 599
600 Barbara Plank. The ‘problem’ of human label variation: On ground truth in data, modeling and
601 evaluation. *arXiv preprint arXiv:2211.02570*, 2022.
- 602
603 Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q Weinberger. Identifying mislabeled
604 data using the area under the margin ranking. In H. Larochelle, M. Ranzato, R. Hadsell, M. F.
605 Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp.
606 17044–17056. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper/2020/file/c6102b3727b2a7d8b1bb6981147081ef-Paper.pdf)
607 [paper/2020/file/c6102b3727b2a7d8b1bb6981147081ef-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/c6102b3727b2a7d8b1bb6981147081ef-Paper.pdf).
- 608 Robert Plutchik. A general psychoevolutionary theory of emotion. In *Theories of emotion*, pp. 3–33.
609 Elsevier, 1980.
- 610
611 Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. Towards empathetic
612 open-domain conversation models: a new benchmark and dataset. In *ACL*, 2019.
- 613
614 Mengye Ren, Wenyan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for
615 robust deep learning. In *International conference on machine learning*, pp. 4334–4343. PMLR,
616 2018.
- 617
618 Mobashir Sadat and Cornelia Caragea. Hierarchical multi-label classification of scientific documents.
619 In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference*
620 *on Empirical Methods in Natural Language Processing*, pp. 8923–8937, Abu Dhabi, United Arab
621 Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
emnlp-main.610. URL <https://aclanthology.org/2022.emnlp-main.610>.
- 622
623 Shreyas Saxena, Oncel Tuzel, and Dennis DeCoste. Data parameters: A new family of parameters for
624 learning a differentiable curriculum. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc,
625 E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Cur-
626 ran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/](https://proceedings.neurips.cc/paper/2019/file/926ffc0ca56636b9e73c565cf994ea5a-Paper.pdf)
627 [file/926ffc0ca56636b9e73c565cf994ea5a-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/926ffc0ca56636b9e73c565cf994ea5a-Paper.pdf).
- 628
629 Klaus R Scherer and Harald G Wallbott. Evidence for universality and cultural variation of differential
630 emotion response patterning. *Journal of personality and social psychology*, 66(2):310, 1994.
- 631
632 Abhinav Shrivastava, Tomasz Malisiewicz, Abhinav Gupta, and Alexei A Efros. Data-driven visual
633 similarity for cross-domain image matching. *ACM Trans. Graph.*, 30(6):154, 2011.
- 634
635 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and
636 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank.
637 In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp.
638 1631–1642, Seattle, Washington, USA, October 2013a. Association for Computational Linguistics.
639 URL <https://aclanthology.org/D13-1170>.
- 640
641 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and
642 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank.
643 In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp.
644 1631–1642, 2013b.
- 645
646 Tiberiu Sosea and Cornelia Caragea. Canceremo: A dataset for fine-grained emotion detection.
647 In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*
(*EMNLP*), pp. 8892–8904, 2020.
- 648
649 Carlo Strapparava, Rada Mihalcea, and Alberto Battocchi. A parallel corpus of music and lyrics
650 annotated with emotions. In *LREC*, pp. 2343–2346. Citeseer, 2012.

- 648 Varsha Suresh and Desmond Ong. Not all negatives are equal: Label-aware contrastive loss for fine-
649 grained text classification. In *Proceedings of the 2021 Conference on Empirical Methods in Natural*
650 *Language Processing*, pp. 4381–4394, Online and Punta Cana, Dominican Republic, November
651 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.359. URL
652 <https://aclanthology.org/2021.emnlp-main.359>.
- 653 Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A.
654 Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with training
655 dynamics. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language*
656 *Processing (EMNLP)*, pp. 9275–9293, Online, November 2020. Association for Computational
657 Linguistics. doi: 10.18653/v1/2020.emnlp-main.746. URL [https://aclanthology.org/](https://aclanthology.org/2020.emnlp-main.746)
658 [2020.emnlp-main.746](https://aclanthology.org/2020.emnlp-main.746).
- 659 Jörg Tiedemann and Santhosh Thottingal. OPUS-MT — Building open translation services for the
660 World. In *Proceedings of the 22nd Annual Conferenec of the European Association for Machine*
661 *Translation (EAMT)*, Lisbon, Portugal, 2020.
- 662 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
663 Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
664 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 665 Wenbo Wang, Lu Chen, Krishnaprasad Thirunarayan, and Amit P Sheth. Harnessing twitter" big
666 data" for automatic emotion identification. In *2012 International Conference on Privacy, Security,*
667 *Risk and Trust and 2012 International Confernece on Social Computing*, pp. 587–592. IEEE, 2012.
- 668 Xinyi Wang, Hieu Pham, Zihang Dai, and Graham Neubig. SwitchOut: an efficient data augmentation
669 algorithm for neural machine translation. In *Proceedings of the 2018 Conference on Empirical*
670 *Methods in Natural Language Processing*, pp. 856–861, Brussels, Belgium, October-November
671 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1100. URL [https:](https://aclanthology.org/D18-1100)
672 [//aclanthology.org/D18-1100](https://aclanthology.org/D18-1100).
- 673 Hongxin Wei, Huiping Zhuang, Renchunzi Xie, Lei Feng, Gang Niu, Bo An, and Yixuan Li.
674 Mitigating memorization of noisy labels by clipping the model prediction. In Andreas Krause,
675 Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett
676 (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of
677 *Proceedings of Machine Learning Research*, pp. 36868–36886. PMLR, 23–29 Jul 2023a. URL
678 <https://proceedings.mlr.press/v202/wei23e.html>.
- 679 Jiaheng Wei, Zhaowei Zhu, Tianyi Luo, Ehsan Amid, Abhishek Kumar, and Yang Liu. To aggregate
680 or not? learning with separate noisy labels. In *Proceedings of the 29th ACM SIGKDD Conference*
681 *on Knowledge Discovery and Data Mining, KDD '23*, pp. 2523–2535, New York, NY, USA, 2023b.
682 Association for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599522.
683 URL <https://doi.org/10.1145/3580305.3599522>.
- 684 Qi Wei, Lei Feng, Haoliang Sun, Ren Wang, Chenhui Guo, and Yilong Yin. Fine-grained classification
685 with noisy labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
686 *Recognition (CVPR)*, pp. 11651–11660, June 2023c.
- 687 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
688 Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick
689 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger,
690 Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural
691 language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural*
692 *Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. Association
693 for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6. URL [https://www.](https://www.aclweb.org/anthology/2020.emnlp-demos.6)
694 [aclweb.org/anthology/2020.emnlp-demos.6](https://www.aclweb.org/anthology/2020.emnlp-demos.6).
- 695 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
696 deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- 697 Mike Zhang and Barbara Plank. Cartography active learning. *arXiv preprint arXiv:2109.04282*,
698 2021.

Dataset	Empathetic Dialogues (wF1)	GoEmotions (wF1)	ISEAR (wF1)	CancerEmo (wF1)	RCV1 (wF1)
BASE	—	—	—	—	—
E2L	—	—	—	—	—
H2L	—	—	—	—	—
AMG	—	—	—	—	—
NSE	—	—	—	—	31.4 ± 1.7
AUM	10.4 ± 0.6	17.5 ± 1.3	27.8 ± 0.8	41.8 ± 0.9	32.5 ± 1.3
LCL	—	—	—	—	—
SCL	—	—	—	—	—
DISC	<u>14.1 ± 1.7</u>	<u>19.6 ± 0.7</u>	31.4 ± 0.5	<u>47.6 ± 0.9</u>	33.7 ± 1.5
UNICON	13.7 ± 1.4	17.4 ± 1.2	33.1 ± 0.9	46.5 ± 0.9	34.6 ± 1.5
LANE	14.6 ± 1.2	20.5 ± 0.9	35.1 ± 0.7	50.1 ± 0.6	38.2 ± 1.7
DATASET	SciHTC (MF1)	SST-5 (Acc)	Amazon Review (Acc)	Yelp (Acc)	Yahoo (Acc)
BASE	—	—	—	—	—
E2L	—	—	—	—	—
H2L	—	—	—	—	—
AMG	—	—	—	—	—
NSE	14.8 ± 1.5	41.6 ± 2.3	—	44.7 ± 2.6	—
AUM	17.2 ± 1.4	42.6 ± 1.5	51.4 ± 1.1	52.6 ± 1.8	42.7 ± 1.9
LCL	—	—	—	—	—
SCL	—	—	—	—	—
DISC	18.5 ± 2.3	43.8 ± 1.8	52.9 ± 1.9	53.8 ± 2.3	44.7 ± 2.1
UNICON	19.6 ± 1.5	43.1 ± 1.6	55.2 ± 1.3	53.9 ± 1.7	44.7 ± 2.1
LANE	20.5 ± 1.5	45.7 ± 1.3	56.8 ± 2.2	56.2 ± 2.3	46.3 ± 2.5

Table 6: Performance of LANE on the the ten benchmark datasets under 40% label noise. The reported results are averaged across five runs and standard deviations are provided. Best results are shown in **bold blue** and second best are underlined. Results marked with — indicate that the model did not converge.

A DATASETS WITH 40% LABEL NOISE

We show in Table 6 results on the 40% noise (40N) datasets. Results marked with - indicate that the model did not converge. We notice that LANE stays effective across the ten datasets, and we observe that AUM yields poor results on this dataset with very high amounts of noise, indicating that it may not work in high-noise setups. For example, AUM outperforms DISC by an average of 1.5% on 20N across the datasets whereas DISC outperforms AUM on 40N by a significant 2.9%. Critically, LANE outperforms both DISC and AUM on 40N by an average of 2.2% and 6.2%, respectively.

B ERROR ANALYSIS

To provide additional insights into our method, we show in Figure 2 a confusion matrix of our LANE approach compared with LANE^{-alm} and a base BERT model on the 20N ISEAR dataset. We make a few observations. First, we note that LANE^{-alm} improves the capabilities of the model over the plain BERT to distinguish between closely related emotions. For example, we see that there are significantly fewer prediction errors confusing disgust and anger or sadness and anger. This result aligns with the purpose of the contrastive loss in LANE^{-alm}, which tries to produce language representations that are useful for distinguishing between confusable classes such as anger, disgust, and sadness. Interestingly, we notice that while the performance on closely confusable classes improves, the performance of the model on opposite or more dissimilar classes degrades. For instance, we observe that the model predicts significantly more examples with disgust as true label in the joy class. However, our LANE solves this drawback and we note that the confusability between opposite classes is considerably improved, outperforming the base BERT model as well substantially. Thus, the combination of contrastive learning with our label-aware approach for learning under label noise is extremely effective, denoting that the two components are complementary by nature: while LANE^{-alm} improves the capabilities of the model of distinguishing between easily confusable classes, our full LANE model improves on both highly confusable/overlapping classes and distant classes.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

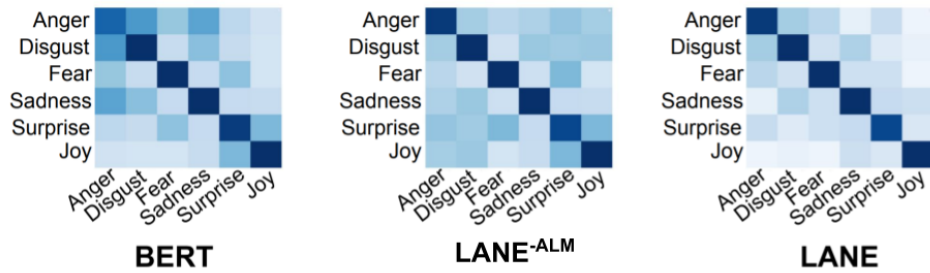


Figure 2: Confusion matrices on the ISEAR dataset created using 20% noise. We compare LANE with a vanilla BERT base model and LANE^{-alm} ablation.