

LiNC: LIGHTWEIGHT NOISE CORRECTION VIA ADAPTIVE LABEL REFINEMENT

Anonymous authors

Paper under double-blind review

ABSTRACT

Medical imaging datasets often contain label noise due to factors such as inter-rater variability, annotation errors, and ambiguous cases, which may severely undermine the reliability and clinical effectiveness of machine learning models trained using those datasets. To address this challenge, we introduce Lightweight Noise Correction (LiNC), which is an intuitive and powerful approach that assigns a trainable trust parameter, α_i , to each individual training sample. Initially initialized to fully trust the observed labels, these parameters adaptively shift trust towards model predictions through a gradient-based optimization process, effectively identifying and reducing the impact of noisy labels by correcting them. After this correction process, usual model training is carried out. Our method requires minimal computational overhead, making it practical for widespread adoption in cases where noise is suspected within a dataset. Extensive evaluations on ten medical imaging datasets from the MedMNISTv2 collection reveal significant improvements in classification accuracy and AUROC across various uniform label noise levels (ranging from 0% to 50%) and robust detection of mislabeled samples, underscoring LiNC’s potential to improve noisy machine learning.

1 INTRODUCTION

Within healthcare, medical imaging is essential in supporting clinical tasks such as diagnosis, treatment planning, and disease monitoring. Recent advancements in deep learning have significantly improved medical image analysis by automating the detection and classification of various medical conditions. However, these advancements heavily depend on the availability of accurately labeled datasets. Machine learning models tend to severely degrade in performance when trained on noisy data. On the other hand, label noise is prevalent in healthcare datasets due to inconsistent annotations, ambiguous findings, and human errors during the annotation process. As a result, these kinds of inaccuracies can pose severe limitations to the reliability of machine learning models and may potentially have adverse effects on patient outcomes if not dealt with properly.

Current approaches to handling label noise predominantly include methods such as sample weighting, output regularization, non-self label correction, and self-label correction, as discussed in Related Work. While some of these methods may be effective at mitigating the effects of label noise, not all of them correct label noise and even when they do, they often may need a high computational load, tedious hyperparameter tuning, additional loss terms, clean validation sets, and may not even take the evolving training dynamics into account. This limits their use in practical clinical settings.

To overcome these limitations, we propose Lightweight Noise Correction (LiNC), a simple yet powerful methodology designed to improve accuracy, interpretability, and reliability in medical imaging classification tasks. LiNC introduces a differentiable, per-sample parameter, α_i , which initially places full trust in the observed labels but dynamically shifts trust towards the model’s predictions based on the learned confidence levels, which is automatically learned through the training process. This approach not only corrects noisy labels effectively but also provides direct insight into label reliability, through the analysis of the α_i parameters.

054 In summary, we introduce a new method, LiNC, that:
055

- 056 1. adds one extra trainable parameter per sample and only needs to update the parameters
057 per-batch using manual gradient update (maintains similar time/space complexity);
- 058 2. performs label correction based on learned knowledge to improve model training and
059 enhance test performance (does not require weighting or pruning valuable data);
- 060 3. does not require a clean validation set or meta-set for the method (may not necessarily be
061 available in the health domain); and,
- 062 4. is lightweight, flexible, and can be easily incorporated into existing training or fine-tuning
063 workflows.
064

066 2 RELATED WORK 067

068 The challenge of training deep neural networks in the presence of label noise has motivated several
069 methods designed to either suppress noise or directly correct errors. Although sample weighting
070 methods (Katharopoulos & Fleuret, 2019; Ren et al., 2019; Xu et al., 2021; Kong et al., 2021; Zhu
071 et al., 2022; Zhou et al., 2023; Wu et al., 2023; Zhou et al., 2023; Wu & Li, 2024; Jain et al., 2024) aim
072 to reduce the weight of certain samples based on difficulty (due to label noise, input noise, ambiguous
073 samples, etc.), label modification can help more directly correct and improve the training trajectory.

074 As mentioned by Wang et al. (2022) in great detail, there are two main categories of label modification:
075 output regularization, which works by adjusting the confidence levels of the targets, and label
076 correction, which involves modifying the labels based on some information.

077 Output regularization aims to prevent models from becoming overly confident in their predictions
078 with methods such as: label smoothing (Müller et al., 2020) and confidence penalty (Pereyra et al.,
079 2017). Label smoothing introduces uncertainty into labels by replacing the hard one-hot target
080 distribution with a soft, convex combination of the original label and a uniform distribution across
081 all classes to help reduce overfitting and encourage better generalization, particularly in clean-label
082 scenarios. Similarly, confidence penalty discourages low-entropy output distributions. Both of these
083 methods help improve calibration and robustness on clean or mildly corrupted datasets. *However,*
084 *they do not perform label correction, so they are not as useful in high-noise scenarios where incorrect*
085 *labels need to be explicitly identified and corrected.*

086 To address this, label correction strategies have been developed to revise the labels based on more
087 reliable estimates of the true label distribution. These approaches can be further divided into non-self
088 and self label correction methods.

089 Non-self label correction relies on external models - called teacher networks - to generate soft labels
090 that are used as improved supervision for the student network. An example of non-self learning
091 correction is knowledge distillation (Hinton et al., 2015), where a pretrained or concurrently trained
092 teacher model provides probabilistic targets for the student model. This uses knowledge from a
093 potentially more stable or better-calibrated model to hopefully remove noise in the observed labels
094 and improve training dynamics. *However, non-self label correction approaches inherently require*
095 *significant overhead and are susceptible to performance degradation if the teacher model is biased*
096 *or miscalibrated.*

097 Self label correction approaches aim to refine labels using the model’s own predictions, removing the
098 need for additional models and improving chances for scalability and practical deployment. These
099 methods are simpler and more end-to-end compatible. One such approach proposed by Lee et al.
100 (2013) replaces the observed label with the model’s current prediction by using the highest-probability
101 class. In some cases (Vyas et al., 2020), soft-labels may make training unstable. *However, this*
102 *method can suffer from confirmation bias, especially in early training phases when model predictions*
103 *are still unstable and inaccurate.*

104 In some of these methods, incorrect predictions may be reinforced in self label correction methods,
105 causing a feedback loop that worsens noise. To mitigate this, bootstrapping methods (Reed et al.,
106 2014) blend the observed label with the model’s current predictions using a convex combination.
107 This interpolation is controlled by a fixed parameter that determines the trust in the model’s predic-
tions. These methods offer a trade-off between retaining potentially noisy labels and incorporating

108 corrections by the model. *However, the fixed nature of the parameter does not take the evolving*
 109 *trustworthiness of the model’s predictions over the course of training into consideration.*

110
 111 Other self label correction strategies attempt to address this by introducing stage-wise learning.
 112 Tanaka et al. (2018); Yuan et al. (2019); Lu et al. (2023) split training into stages where the model
 113 first learns from the observed, potentially noisy labels and later on, the model’s predictions are fully
 114 trusted and used as labels. These models effectively deal with avoiding the need for a fixed trust
 115 parameter. *However, they still suffer from significant overhead.*

116 Most of these methods require additional models, hyperparameter tuning, manually chosen scores,
 117 trust in a potentially untrustworthy model or potentially unstable stage of training, and/or a higher
 118 computational load, making these methods less scalable and harder to deploy in real-world scenarios,
 119 especially in high-stakes domains like healthcare and medical imaging.

120 Therefore, in our method, we introduce a single extra differentiable parameter for each training
 121 sample that will use concepts from self label correction to either give more weight to the observed
 122 label or the model prediction after a few warmup epochs without requiring much additional overhead.
 123 Then, the method uses Otsu thresholding Otsu et al. (1975) on these smooth/convex parameters to
 124 separate samples into ones that will stick to their original label and ones that will switch to the label
 125 provided by the model. Unlike prior work, our method does not rely on extensive tuning, auxiliary
 126 models, or discrete training phases, and can be seamlessly integrated into standard training loops
 127 after a brief warmup phase, a brief correction phase, and then a final training phase offering a scalable
 128 and model-agnostic solution to label noise correction. In addition, our method also has potential to
 129 be very strong with pre-trained models when there is a need for fine-tuning rather than training from
 130 scratch for the best results.

131 3 METHOD

132 In this section, we provide a description of our method along with our overall training procedure. Let
 133 us consider the following supervised multi-class classification problem.

134 Let $\mathcal{D} = \{(x_i, \tilde{y}_i)\}_{i=1}^N$ be the training set, where (x_i, \tilde{y}_i) are the inputs and potentially noisy, observed
 135 targets. In our case, $x_i \in \mathcal{X}$, where \mathcal{X} is the input space of images and $\tilde{y}_i \in \mathcal{Y} = \{0, \dots, c - 1\}$, is
 136 the output space with $c \in \mathbb{N}$ such that $c \geq 2$ is the number of classes. We also have the parameters α_i
 137 for $i \in \{1, \dots, N\}$, i.e. one for each training sample. Let $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ be the neural network model
 138 and θ be its parameters.

139 During training, our method optimizes both the network parameters and the α_i parameters. The
 140 training process is detailed in Algorithm 1.

141 The core mechanism underlying our noise correction framework is the interpolation of training targets
 142 using a per-sample trust parameter $\alpha_i \in [0, 1]$, which dynamically adjusts the degree of trust placed
 143 on the observed label versus the model’s own prediction:

$$144 \tilde{y}_i = (1 - \alpha_i) \cdot s_i + \alpha_i \cdot \tilde{y}_i$$

145 When $\alpha_i = 1$: the model fully trusts the observed label \tilde{y}_i .

146 When $\alpha_i = 0$: the model completely trusts its own prediction s_i .

147 Here, \tilde{y}_i represents the original observed label, and s_i denotes the model’s predicted probability
 148 distribution for a given sample. By treating α_i as a learnable trust coefficient that varies across
 149 samples and over time, we enable the model to self-modulate how much it sticks to the observed
 150 label versus its internal belief. This dynamic trust mechanism becomes useful in the presence of label
 151 noise when the model is given flexibility to assign more importance to its own prediction if it detects
 152 conflicts with the provided label.

153 During each training iteration, the α_i values are updated through a manual gradient descent step (as
 154 shown in Line 19 of Algorithm 1), with the objective of refining label reliability. We enforce a discrete
 155 decision and maintain stability during the training process (avoiding meaningless α_i values), after a

Algorithm 1 Training with LiNC

```

162 1: Inputs: training data  $\mathcal{D}$ , classifier model  $f_\theta$ , trust parameters  $\alpha_i$ , number of warmup epochs  $w$ ,
163 2:     optimizer for all  $\alpha_i$  and  $\theta$ , learning rate  $\alpha_{lr}$ , weight decay  $\alpha_{wd}$ 
164 3: Output: classifier model trained using LiNC
165 4: for epoch  $\leftarrow 1$  to  $n$  do
166 5:   for batch  $(x, \tilde{y}) \leftarrow \text{dataloader}(\mathcal{D})$  do
167 6:     set requires_grad to True on  $\theta$  and  $\{\alpha_j\}$  for current batch  $j \in \{1, \dots, N\}$ 
168 7:     let  $x_j = x$  and  $\tilde{y}_j = \tilde{y}$ 
169 8:      $s_j = \text{softmax}(f_\theta(x_j))$ 
170 9:     if epoch  $< w$  then
171 10:        $\tilde{y}_j = (1 - \{\alpha_j\}) * s_j + (\{\alpha_j\}) * \tilde{y}_j$ 
172 11:     end if
173 12:     if epoch ==  $w$  then
174 13:        $\tilde{y}_j = s_j$ , if  $\{\alpha_j\} < \text{otsu}(\alpha_i)$ 
175 14:     end if
176 15:     calculate train loss  $\mathcal{L} = \ell(f_\theta(x), \tilde{y})$ 
177 16:      $\mathcal{L}.\text{backward}()$ 
178 17:     optimizer.step() to update  $\theta$ 
179 18:     if epoch  $< w$  then
180 19:        $\{\alpha_j\} = \{\alpha_j\} - \alpha_{lr} * \nabla_{\{\alpha_j\}} \mathcal{L}$ 
181 20:        $\{\alpha_j\} = \{0 \text{ if } \alpha_j < 0 \text{ and } 1 \text{ if } \alpha_j > 1\}$ 
182 21:     end if
183 22:     optimizer.zero_grad()
184 23:   end for
185 24: end for
186 25: return  $f_\theta(x)$ 

```

few warmup epochs. We project the updated α_i values back to $[0, 1]$ in Line 20. Noise correction happens once in Line 13 after the warmup epochs.

Also, the behavior of the α_i parameters evolves throughout training. Samples with noisy or inconsistent labels tend to rapidly shift their α_i values towards 0, whereas clean samples generally stabilize with α_i values near 1. This separation implicitly segments the dataset into reliable and unreliable subsets, without requiring any prior knowledge of which samples are mislabeled. Furthermore, because the α_i updates are localized and differentiable, they can be used easily in any model training with minimal effort.

Why Otsu on α ? We threshold α_i using Otsu’s method to obtain a data-adaptive, parameter-free split between trusting the observed labels and trusting the model prediction. Otsu exploits bimodality that naturally emerges as clean samples keep $\alpha_i \approx 1$ while mislabels move toward 0. In addition, Otsu is much faster than most alternatives, differentiable, and introduces no extra hyperparameters.

4 EXPERIMENTAL SETUP

Our experiments are conducted on ten 2D datasets with different medical imaging modalities focusing on multi-class classification tasks: PathMNIST (Kather et al., 2019), DermaMNIST (Tschandl et al., 2018; Codella et al., 2019), OCTMNIST (Keremany et al., 2018), PneumoniaMNIST (Keremany et al., 2018), BreastMNIST (Al-Dhabyani et al., 2020), BloodMNIST (Acevedo et al., 2020), TissueMNIST (Ljosa et al., 2012), OrganAMNIST (Bilic et al., 2023; Xu et al., 2019), OrganCMNIST (Bilic et al., 2023; Xu et al., 2019), and OrganSMNIST (Bilic et al., 2023; Xu et al., 2019), from the MedMNISTv2 collection (Yang et al., 2021; 2023; Doerrich et al., 2024). All the images are of size 224×224 and we do not use any augmentations.

Using ImageNet-21K pretraining, we trained on 224×224 MedMNISTv2 inputs for 100 epochs with a batch size of 128, learning rate of $1e-4$, and no weight decay. We used an Adam optimizer (Kingma, 2014) and a multi-step learning rate scheduler was applied, decaying the rate by a factor of 0.1 at epochs 50 and 75, following the setup in Yang et al. (2023).

We used a warmup period of 10 epochs to ensure that the model briefly learns about the mislabels in the data before using LiNC. We initially set the α_i parameters all to 1 (fully trust the observed labels), with learning rate $\alpha_{lr} = 1$ and weight decay $\alpha_{wd} = 1e - 3$ for manual gradient descent. Our method is not too sensitive to this learning rate choice as the learning rate needs to be large enough to push the α_i for mislabels to 0 since generally, the α_i for correct labels do not change much at all. We inject various levels of uniform noise (0%, 10%, 20%, 30%, 40%, 50%) to study the effectiveness of LiNC. We analyze the classification accuracy and AUROC with and without LiNC and study its ability to identify mislabels compared to several baselines.

5 RESULTS

Method	AUROC	
	epoch 2	epoch 10
CNLCU-S	0.6688	0.6829
GraND	0.6109	0.5785
Data-IQ	0.5275	0.6763
DataMaps	0.5000	0.4084
EL2N	0.6688	0.5606
AUM	0.6364	0.6794
α_i	0.7923	0.8012

Table 1: We report the AUROC scores of various methods in being able to identify mislabels when compared to the α_i parameters from our method after the warmup epochs on the BreastMNIST dataset with 50% uniform label noise.

In Table 1, we compare the α_i parameters from our method to several other baseline methods that have been effective in identifying noisy or otherwise “difficult” data: AUM (Pleiss et al., 2020), DataMaps (Swayamdipta et al., 2020), Data-IQ (Seedat et al., 2022), EL2N (Paul et al., 2021), GraND (Paul et al., 2021), and CNLCU-S (Xia et al., 2021). Note that most of these methods also require quite a bit more work to calculate and/or maintain whereas the α_i parameters from our method are simply maintained and updated during the training process.

Our method achieves the highest AUCROC score of 0.7923 after epoch 2 and 0.8012 after epoch 10, substantially outperforming all baselines. This result indicates that the learned α_i values are highly predictive of label correctness, providing strong evidence that our adaptive label correction mechanism is effective at isolating noisy labels early in training (just after 10 warmup epochs). These results highlight the practical benefits of our proposed method, which identifies trustworthiness through the α_i parameters without requiring auxiliary models or extensive post-processing. With this, we have shown that label correction can be formulated as a differentiable, per-sample trust mechanism which effectively learns to identify noisy supervision from within the training loop itself.

We observe that the per-sample trust parameters α_i , exhibit a highly discriminative relationship between clean and noisy labels. Specifically, values of α_i are driven toward 1 for samples where the observed label is likely to be correct and values approach 0 for samples where the model prediction is deemed more reliable, indicating a lower degree of trust in the observed label. This happens consistently even under the use of a relatively large learning rate of 1 and weight decay of $1e - 3$ for the α_i parameters. Despite the aggressive gradient updates, the optimization process maintains stability and convergence toward meaningful values that reflect the underlying quality of each label.

Therefore, our sample-wise label correction mechanism can dynamically infer which instances are mislabeled and adjust supervision accordingly. Even under minimal constraints and with highly localized updates, per-sample trust parameters can meaningfully disentangle noisy and clean samples and improve robustness to noise during training.

Table 2 shows the accuracy results of running the model with and without LiNC across ten 2D medical imaging datasets, of varying modalities, from the MedMNISTv2 (Yang et al., 2021; 2023; Doerrich et al., 2024) collection under varying degrees of uniform label noise (0%-50%). For each dataset, we report the baseline accuracy achieved without LiNC, and the accuracy improvements obtained

Dataset	Acc. (%)					
	0% Noise	10% Noise	20% Noise	30% Noise	40% Noise	50% Noise
PathMNIST	95.64	94.96 ↑ 0.48	94.70 ↑ 0.20	94.19 ↑ 0.41	93.27 ↑ 0.18	92.82 ↑ 0.92
DermaMNIST	73.70	72.01 ↑ 2.37	72.60 ↑ 1.98	70.87 ↑ 2.63	70.50 ↑ 1.37	69.87 ↑ 1.83
OCTMNIST	75.29	70.69 ↓ 0.15	69.84 ↑ 1.22	72.54 ↓ 3.10	68.08 ↓ 2.73	65.06 ↓ 0.39
PneumoniaMNIST	88.13	87.57 ↑ 3.10	87.14 ↑ 1.77	85.83 ↑ 3.05	86.03 ↓ 2.70	86.99 ↑ 0.73
BreastMNIST	83.26	84.26 ↑ 1.17	84.65 ↑ 0.39	80.92 ↓ 0.40	72.82 ↑ 9.66	69.87 ↑ 8.26
BloodMNIST	98.04	97.32 ↑ 0.78	96.21 ↑ 1.33	95.48 ↑ 1.60	94.15 ↑ 1.90	92.78 ↑ 1.78
TissueMNIST	67.16	65.95 ↓ 1.42	64.18 ↓ 1.03	63.31 ↓ 1.04	61.83 ↓ 0.76	60.01 ↓ 0.12
OrganAMNIST	94.37	92.78 ↑ 0.77	91.77 ↑ 0.60	90.64 ↑ 1.08	90.49 ↑ 0.66	89.27 ↑ 0.24
OrganCMNIST	87.44	85.12 ↑ 2.20	84.27 ↑ 2.49	82.22 ↑ 2.64	80.76 ↑ 2.55	78.61 ↑ 2.97
OrganSMNIST	77.49	75.43 ↑ 1.61	73.65 ↑ 1.92	71.52 ↑ 2.92	70.21 ↑ 2.97	65.10 ↑ 3.96

Table 2: Accuracy on ten 2D datasets from MedMNISTv2 with different levels of uniform noise.

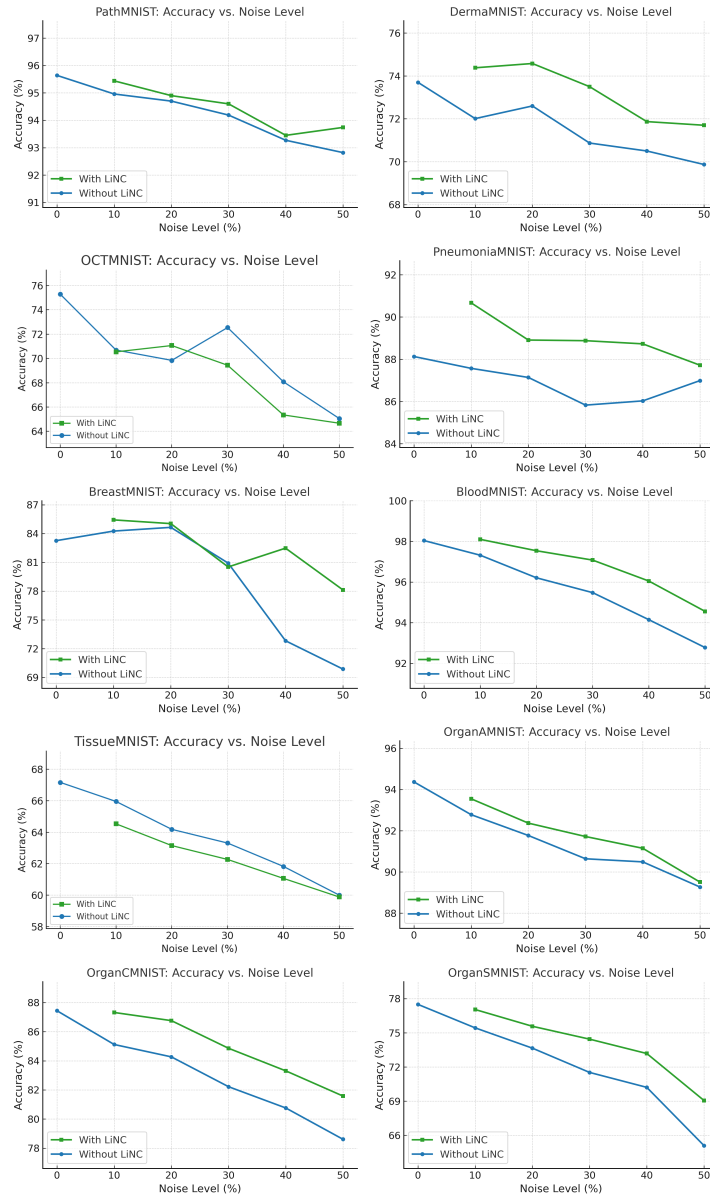


Figure 1: Accuracy with and without LiNC on ten 2D datasets from MedMNISTv2.

Dataset	AUROC					
	0% Noise	10% Noise	20% Noise	30% Noise	40% Noise	50% Noise
PathMNIST	0.9968	0.9934 \uparrow 0.0033	0.9912 \uparrow 0.0050	0.9898 \uparrow 0.0053	0.9881 \uparrow 0.0052	0.9862 \uparrow 0.0057
DermaMNIST	0.8936	0.8630 \downarrow 0.0063	0.8536 \uparrow 0.0065	0.8338 \uparrow 0.0278	0.8004 \uparrow 0.0489	0.7845 \uparrow 0.0395
OCTMNIST	0.9817	0.9691 \uparrow 0.0139	0.9626 \uparrow 0.0127	0.9580 \uparrow 0.0114	0.9569 \downarrow 0.0019	0.9380 \uparrow 0.0089
PneumoniaMNIST	0.9646	0.9706 \uparrow 0.0092	0.9659 \uparrow 0.0063	0.9666 \uparrow 0.0056	0.9516 \uparrow 0.0164	0.9441 \uparrow 0.0146
BreastMNIST	0.8678	0.8680 \downarrow 0.0046	0.8594 \uparrow 0.0001	0.8446 \downarrow 0.0173	0.7878 \uparrow 0.0370	0.7621 \downarrow 0.0023
BloodMNIST	0.9984	0.9976 \uparrow 0.0012	0.9967 \uparrow 0.0016	0.9954 \uparrow 0.0019	0.9933 \uparrow 0.0023	0.9908 \uparrow 0.0020
TissueMNIST	0.9074	0.9035 \downarrow 0.0363	0.8942 \downarrow 0.0399	0.8877 \downarrow 0.0426	0.8777 \downarrow 0.0486	0.8623 \downarrow 0.0403
OrganAMNIST	0.9865	0.9940 \uparrow 0.0015	0.9929 \uparrow 0.0018	0.9926 \uparrow 0.0013	0.9912 \uparrow 0.0006	0.9889 \uparrow 0.0011
OrganCMNIST	0.9928	0.9833 \uparrow 0.0032	0.9819 \uparrow 0.0038	0.9783 \uparrow 0.0044	0.9756 \uparrow 0.0037	0.9723 \uparrow 0.0043
OrganSMNIST	0.9670	0.9635 \uparrow 0.0049	0.9564 \uparrow 0.0070	0.9529 \uparrow 0.0083	0.9496 \uparrow 0.0087	0.9387 \uparrow 0.0092

Table 3: AUROC on ten 2D datasets from MedMNISTv2 with different levels of uniform noise.

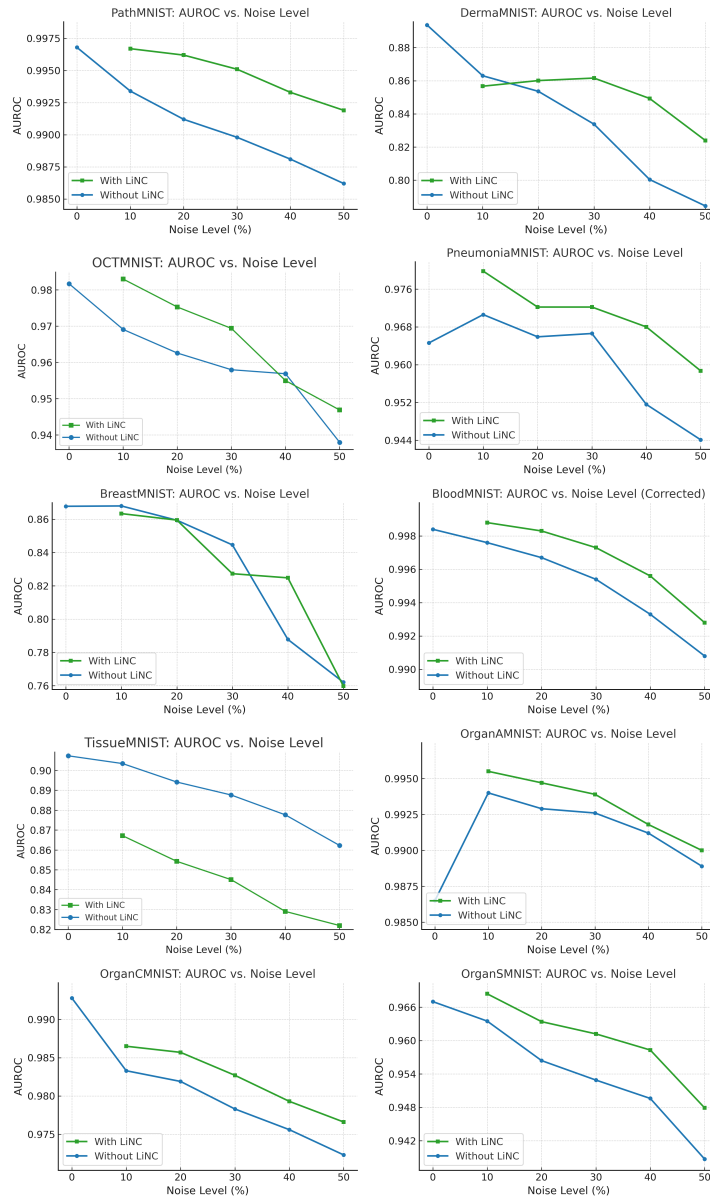


Figure 2: AUROC with and without LiNC on ten 2D datasets from MedMNISTv2.

when applying LiNC. Overall, the results consistently demonstrate that LiNC robustly improves classification accuracy across all datasets and noise levels examined. Specifically, the method

378 significantly mitigates the negative impact of label noise, with accuracy improvements becoming
379 more pronounced as the noise level increases. This finding underscores LiNC’s inherent capability
380 of identifying mislabels and dynamically enabling the network to rely more strongly on its own
381 progressively confident predictions rather than noisy training labels. It is especially noteworthy that
382 in almost all cases where there is injected noise, the improvement in accuracy with LiNC is similar to
383 the performance without LiNC when there is no injected noise. Figure 1 is a visual representation of
384 Table 2. You can see that in almost all cases (except with OCTMNIST and TissueMNIST, which
385 seem to be the hardest tasks on the MedMNISTv2 dataset), there is a significant increase in accuracy
386 after correcting noise with LiNC.

387 Table 3 reports AUROC scores across the datasets. We see that large-scale or relatively easier datasets
388 remain highly robust, retaining AUROC scores above 0.98 even at 50% noise, with only marginal
389 fluctuations. More noise-sensitive datasets with AUROC scores less than 0.80 show degradations as
390 noise increases, with TissueMNIST dropping from 0.9074 at 0% noise to 0.8623 at 50% while using
391 LiNC. The remaining datasets show small but consistent improvements. It is also important to note
392 that AUROC does not uniformly decrease with increasing noise. Figure 2 is a visual representation of
393 Table 3.

394 6 DISCUSSION

395 Our results demonstrate the effectiveness of LiNC in mitigating label noise across various medical
396 imaging datasets and different noise levels. The method consistently improved accuracy, indicating
397 its robustness and adaptability in diverse medical imaging modalities, from dermatology to pathology.
398 Notably, LiNC exhibited remarkable performance improvements even under extreme noise conditions
399 (up to 50%), demonstrating its strength in accurately identifying and correcting mislabeled samples,
400 which is crucial in high-stakes healthcare scenarios.

401 LiNC achieves these benefits with minimal computational overhead. By introducing only a single
402 additional trainable parameter per sample, which dynamically adjusts based on evolving training
403 dynamics, LiNC provides a scalable solution applicable to real-world clinical settings where compu-
404 tational resources and model interpretability are critical considerations. Model developers can use
405 the α_i parameters as an easy way to perform human-in-the-loop audits of health data. Furthermore,
406 LiNC’s simplicity and flexibility mean it can be integrated into existing training workflows, including
407 scenarios involving fine-tuning pretrained models or foundation models.

408 The adaptive nature of the parameters α_i allows LiNC to respond dynamically to the training progress.
409 Our analysis of α_i indicates its effectiveness in differentiating between noisy and clean labels, as
410 evidenced by consistently high AUCROC scores across different noise levels. This adaptability
411 not only enhances accuracy but also provides interpretability into label reliability, offering valuable
412 insights into dataset quality and annotation trustworthiness.

413 One potential limitation of LiNC is its reliance on self label correction. There may be cases where
414 particularly challenging datasets or poorly initialized models could affect early-stage predictions,
415 potentially impacting corrections. Future work could explore adaptive warmup strategies or more
416 complex initializations to further enhance reliability.

417 Since LiNC significantly reduces the negative impact of noisy labels, we aim to investigate hybrid
418 approaches that combine LiNC with active learning and human-in-the-loop verification to study
419 dataset quality and maximize clinical applicability.

420 7 CONCLUSION

421 We presented LiNC, a lightweight, adaptive, and effective approach to address the issue of label noise
422 in medical imaging datasets. By dynamically adjusting the trust placed in observed labels versus
423 model predictions for each training sample, LiNC robustly improves classification accuracy, even in
424 scenarios with substantial label noise. Our extensive experiments across diverse medical datasets
425 demonstrate that LiNC not only enhances accuracy significantly but also provides interpretable
426 insights into dataset reliability. Given its computational efficiency, ease of implementation, and
427 proven effectiveness, LiNC represents a practical and powerful tool for improving the reliability of
428 machine learning models in clinical settings by fixing noisy labels.

432 Future work can explore whether pseudo-labels, in the context of self-supervised or semi-supervised
433 learning, can be fixed using LiNC. Additionally, LiNC may be useful for active learning techniques
434 to identify which samples should be labeled (or relabeled) to most benefit the model.
435

436 REFERENCES 437

- 438 Andrea Acevedo, Anna Merino González, Edwin Santiago Alférez Baquero, Ángel Molina Borrás,
439 Laura Boldú Nebot, and José Rodellar Benedé. A dataset of microscopic peripheral blood cell
440 images for development of automatic recognition systems. *Data in brief*, 30(article 105474), 2020.
441
- 442 Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. Dataset of breast
443 ultrasound images. *Data in brief*, 28:104863, 2020.
- 444 Patrick Bilic, Patrick Christ, Hongwei Bran Li, Eugene Vorontsov, Avi Ben-Cohen, Georgios Kaissis,
445 Adi Szeskin, Colin Jacobs, Gabriel Efrain Humpire Mamani, Gabriel Chartrand, et al. The liver
446 tumor segmentation benchmark (lits). *Medical Image Analysis*, 84:102680, 2023.
- 447 Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David
448 Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion
449 analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging
450 collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019.
451
- 452 Sebastian Doerrich, Francesco Di Salvo, Julius Brockmann, and Christian Ledig. Rethinking model
453 prototyping through the medmnist+ dataset collection. *arXiv preprint arXiv:2404.15786*, 2024.
- 454 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015.
455 URL <https://arxiv.org/abs/1503.02531>.
456
- 457 Nishant Jain, Arun S. Suggala, and Pradeep Shenoy. Improving generalization via meta-learning on
458 hard samples, 2024. URL <http://arxiv.org/abs/2403.12236>.
- 459 Angelos Katharopoulos and François Fleuret. Not all samples are created equal: Deep learning with
460 importance sampling, 2019. URL <http://arxiv.org/abs/1803.00942>.
461
- 462 Jakob Nikolas Kather, Johannes Krisam, Pornpimol Charoentong, Tom Luedde, Esther Herpel, Cleo-
463 Aron Weis, Timo Gaiser, Alexander Marx, Nektarios A Valous, Dyke Ferber, et al. Predicting
464 survival from colorectal cancer histology slides using deep learning: A retrospective multicenter
465 study. *PLoS medicine*, 16(1):e1002730, 2019.
- 466 Daniel S Kermany, Michael Goldbaum, Wenjia Cai, Carolina CS Valentim, Huiying Liang, Sally L
467 Baxter, Alex McKeown, Ge Yang, Xiaokang Wu, Fangbing Yan, et al. Identifying medical
468 diagnoses and treatable diseases by image-based deep learning. *cell*, 172(5):1122–1131, 2018.
469
- 470 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
471 2014.
- 472 Yajing Kong, Liu Liu, Jun Wang, and Dacheng Tao. Adaptive curriculum learning. In *2021*
473 *IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 5047–5056, 2021. doi:
474 10.1109/ICCV48922.2021.00502. URL [https://ieeexplore.ieee.org/document/](https://ieeexplore.ieee.org/document/9709930)
475 [9709930](https://ieeexplore.ieee.org/document/9709930). ISSN: 2380-7504.
- 476 Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for
477 deep neural networks. 2013.
478
- 479 Vebjorn Ljosa, Katherine L Sokolnicki, and Anne E Carpenter. Annotated high-throughput mi-
480 croscopy image sets for validation. *Nature methods*, 9(7):637–637, 2012.
- 481 Peng Lu, Ahmad Rashid, Ivan Kobzyev, Mehdi Rezagholizadeh, and Philippe Langlais. Labo:
482 Towards learning optimal label regularization via bi-level optimization, 2023. URL <https://arxiv.org/abs/2305.04971>.
483
484
- 485 Rafael Müller, Simon Kornblith, and Geoffrey Hinton. When does label smoothing help?, 2020. URL
<http://arxiv.org/abs/1906.02629>.

- 486 Nobuyuki Otsu et al. A threshold selection method from gray-level histograms. *Automatica*, 11
487 (285-296):23–27, 1975.
488
- 489 Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding
490 important examples early in training. *Advances in Neural Information Processing Systems*, 34:
491 20596–20607, 2021.
- 492 Gabriel Pereyra, George Tucker, Jan Chorowski, Lukasz Kaiser, and Geoffrey E. Hinton. Regularizing
493 neural networks by penalizing confident output distributions. *CoRR*, abs/1701.06548, 2017. URL
494 <http://arxiv.org/abs/1701.06548>.
495
- 496 Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q Weinberger. Identifying mislabeled data us-
497 ing the area under the margin ranking. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and
498 H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 17044–17056.
499 Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_](https://proceedings.neurips.cc/paper_files/paper/2020/file/c6102b3727b2a7d8b1bb6981147081ef-Paper.pdf)
500 [files/paper/2020/file/c6102b3727b2a7d8b1bb6981147081ef-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/c6102b3727b2a7d8b1bb6981147081ef-Paper.pdf).
- 501 Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew
502 Rabinovich. Training deep neural networks on noisy labels with bootstrapping. *arXiv preprint*
503 *arXiv:1412.6596*, 2014.
504
- 505 Mengye Ren, Wenyan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for
506 robust deep learning, 2019. URL <http://arxiv.org/abs/1803.09050>.
- 507 Nabeel Seedat, Jonathan Crabbé, Ioana Bica, and Mihaela van der Schaar. Data-IQ: Characterizing
508 subgroups with heterogeneous outcomes in tabular data, 2022. URL [http://arxiv.org/](http://arxiv.org/abs/2210.13043)
509 [abs/2210.13043](http://arxiv.org/abs/2210.13043).
- 510 Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A.
511 Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with training
512 dynamics, 2020. URL <http://arxiv.org/abs/2009.10795>.
- 513 Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. Joint optimization framework
514 for learning with noisy labels. *CoRR*, abs/1803.11364, 2018. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1803.11364)
515 [1803.11364](http://arxiv.org/abs/1803.11364).
- 516 Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of
517 multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9,
518 2018.
- 519 Nidhi Vyas, Shreyas Saxena, and Thomas Voice. Learning soft labels via meta learning. *CoRR*,
520 abs/2009.09496, 2020. URL <https://arxiv.org/abs/2009.09496>.
- 521 Xinshao Wang, Yang Hua, Elyor Kodirov, Sankha Subhra Mukherjee, David A. Clifton, and Neil M.
522 Robertson. Proselflc: Progressive self label correction towards a low-temperature entropy state,
523 2022. URL <https://arxiv.org/abs/2207.00118>.
- 524 Ou Wu and Mengyang Li. Revisiting the effective number theory for imbalanced learning. 36
525 (8):4192–4206, 2024. ISSN 1558-2191. doi: 10.1109/TKDE.2024.3367949. URL <https://ieeexplore.ieee.org/document/10440480>. Conference Name: IEEE Transactions
526 on Knowledge and Data Engineering.
527
- 528 Yinjun Wu, Adam Stein, Jacob Gardner, and Mayur Naik. Learning to select pivotal samples for
529 meta re-weighting, 2023. URL <http://arxiv.org/abs/2302.04418>.
- 530 Xiaobo Xia, Tongliang Liu, Bo Han, Mingming Gong, Jun Yu, Gang Niu, and Masashi Sugiyama.
531 Sample selection with uncertainty of losses for learning with noisy labels, 2021. URL <http://arxiv.org/abs/2106.00445>.
- 532 Da Xu, Yuting Ye, and Chuanwei Ruan. Understanding the role of importance weighting for deep
533 learning, 2021. URL <http://arxiv.org/abs/2103.15209>.

540 Xuanang Xu, Fugen Zhou, Bo Liu, Dongshan Fu, and Xiangzhi Bai. Efficient multiple organ
541 localization in ct image using 3d region proposal network. *IEEE transactions on medical imaging*,
542 38(8):1885–1898, 2019.

543
544 Jiancheng Yang, Rui Shi, and Bingbing Ni. Medmnist classification decathlon: A lightweight automl
545 benchmark for medical image analysis. In *IEEE 18th International Symposium on Biomedical
546 Imaging (ISBI)*, pp. 191–195, 2021.

547 Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and
548 Bingbing Ni. Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image
549 classification. *Scientific Data*, 10(1):41, 2023.

550 Li Yuan, Francis E. H. Tay, Guilin Li, Tao Wang, and Jiashi Feng. Revisit knowledge distillation:
551 a teacher-free framework. *CoRR*, abs/1909.11723, 2019. URL [http://arxiv.org/abs/
552 1909.11723](http://arxiv.org/abs/1909.11723).

553
554 Xiaoling Zhou, Ou Wu, Weiyao Zhu, and Ziyang Liang. Understanding difficulty-based sample
555 weighting with a universal difficulty measure, 2023. URL [http://arxiv.org/abs/2301.
556 04850](http://arxiv.org/abs/2301.04850).

557 Weiyao Zhu, Ou Wu, Fengguang Su, and Yingjun Deng. Exploring the learning difficulty of data
558 theory and measure, 2022. URL <http://arxiv.org/abs/2205.07427>.

559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593