



## Abstract

35 FixMatch is a semi-supervised learning method, which achieves comparable results  
36 with fully supervised learning by leveraging a limited number of labeled data  
37 (pseudo labelling technique) and taking a good use of the unlabeled data (consistency regularization). In this work, we reimplement FixMatch and achieve reasonably comparable performance with the official implementation, which supports that FixMatch outperforms semi-supervised learning benchmarks and demonstrates that the author’s choices with respect to those ablations were experimentally sound. Next, we investigate the existence of a major problem of FixMatch, *confirmation errors*, by reconstructing the batch structure during the training process. It reveals existing confirmation errors, especially the ones caused by *asymmetric noise* in pseudo labels. To deal with the problem, we apply equal-frequency and confidence entropy regularization to the labeled data and add them in the loss function. Our experimental results on CIFAR-10 show that using either of the entropy regularization in the loss function can reduce the asymmetric noise in pseudo labels and improve the performance of FixMatch in the presence of (pseudo) labels containing (asymmetric) noise. Our code is available at the url: <https://github.com/CeliAli/FixMatch>.

## 52 1 Introduction

53 Ghahramani [2020] summarized the reasons for the success of deep learning in his talk given as the  
54 chief scientist in Uber. Firstly, with the availability of large datasets, large models can work well.  
55 Secondly, training such large models with stochastic descent works surprisingly well. Moreover,  
56 staying close to identity (such as ReLU) makes it stable to be trained. The automate differentiation  
57 and a large number of open source softwares make it scale well. Therefore, we can see deep learning  
58 in many applications, such as computer vision, natural language processing, bioinformatics, etc.

59 However, it is not always the case where a huge number of labeled data are available. In some  
60 areas, it is difficult, expensive, or even impossible to have a large labeled dataset, such as medical  
61 images [Kuznetsova et al., 2018]. In this case, it can be difficult to train a Deep Neural Network  
62 (DNN) from scratch with the limited labeled data. Luckily, Tajbakhsh et al. [2016] shows that a  
63 DNN trained based on a pre-trained DNN, fine-tuning, can outperform the one trained from scratch.  
64 Moreover, Semi-Supervised Learning (SSL) is also a common method to deal with the scarcity and  
65 often high acquisition cost of labelled data [von Kügelgen et al., 2020]. SSL efficiently leverages  
66 labeled data and the relation with unlabeled data to train a DNN. Among SSL methods, there is a  
67 class of "match"-based methods, such as FixMatch [Sohn et al., 2020], MixMatch [Berthelot et al.,  
68 2019], ReMixMatch [Kurakin et al., 2020] and DivideMatch [Li et al., 2020]. These methods utilize  
69 the consistency regularization, pseudo-labelling and ensembling methods to boost the performance  
70 with the use of unlabeled data. In fact, they are leveraging prior knowledge to regularize the training  
71 of DNNs. In this project, we focus on reproducing and investigating one of such methods, FixMatch  
72 [Sohn et al., 2020].

73 Nevertheless, SSL is still facing many challenges in theory and in practice. Ben-David et al. show that  
74 “as long as one does not make any assumptions about the behavior of the labels, SSL cannot help much  
75 over algorithms that ignore the unlabeled data.” Moreover, SSL can actually degrade performance if  
76 certain assumptions are not met [Chapelle et al., 2010]. In this line of works, Schölkopf et al. [2012]  
77 consider the problem from a causal modeling perspective and conclude that in fact SSL is impossible  
78 when predicting a target variable from its causes (causal learning) but possible from anti-causal  
79 learning. Recently, the relation of causality and semi-supervised learning is further explored in [von  
80 Kügelgen et al., 2020], i.e., predicting a target variable from both causes and effects at the same  
81 time. Moreover, in the light of consistency regularization and pseudo-labelling, a significant issue of  
82 the "Match"-based methods is *confirmation error*. It happens especially when noisy samples are in  
83 the labeled set. A DNN can keep having lower loss by fitting the noise and be further maintained  
84 after training with the wrong pseudo labels of unlabeled data, which keeps the errors in the model  
85 and limits its generalization and performance [Tarvainen and Valpola]. This problem becomes more  
86 serious in the presence of asymmetric noise in the training labels, which roughly speaking tends to  
87 label a class of data as another specific class.

88 Therefore, in this work, we are not only reimplementing FixMatch, but also investigating whether the  
 89 pseudo labels made by the DNN contain harmful noise leading to confirmation errors. First, we  
 90 design a stable and reliable method to examine the existence of confirmation errors and noisy pseudo  
 91 labels by reconstructing the batch structure. Secondly, we find methods to deal with (asymmetric)  
 92 noise in (pseudo) labels of the training dataset. We reconstruct the batch structure and add an  
 93 equal-frequency entropy regularization on labeled data to the loss function of FixMatch. Moreover,  
 94 we also use a confidence entropy regularization on labeled data to avoid the over-confident prediction.  
 95 It turns out that both entropy regularization is helpful for dealing with the noisy (pseudo) labels (even  
 96 for the asymmetric noise) and confirmation errors. Our experimental results show that

- 97 1. our implementation can achieve almost the same performance even better for low-label  
 98 regimes.
- 100 2. there exists asymmetric noise in the pseudo labels leading to confirmation errors. With such  
 101 pseudo labels, the model is biased which in turn leads to more asymmetric noise in pseudo  
 102 labels.
- 103 3. FixMatch with equal-frequency entropy regularization and FixMatch with confidence en-  
 104 tropy regularization can reduce (asymmetric) noise in the pseudo labels and perform better  
 105 than the baseline FixMatch in the presence of asymmetric noise in (pseudo) labels .

## 106 2 Related work

107 As introduced in Sec. 1, confirmation error is a serious issue of "Match"-based SSL methods and our  
 108 study is mainly about the confirmation error and FixMatch in the presence of noisy (pseudo) labels.  
 109 Therefore, here we mainly introduce the noisy labeling and some related works for dealing with the  
 110 noisy label and confirmation error in SSL.

111 **Noisy labeling and noise-robust loss.** Suppose a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  where  $y_i$  is given by  
 112 noisy labeling. To model noisy labeling process, we have  $p(y_i | \tilde{y}_i)$  where  $\tilde{y}_i$  is the ground truth label  
 113 under the assumption that the noise label is conditionally independent from the input data given the  
 114 ground-true label; formally,  $p(y_i = k | x_i, \tilde{y}_i = j) = p(y_i = k | \tilde{y}_i = j) = \eta_{kj}$ . In general, such noise  
 115 is called class dependent, which is also named as the asymmetric noise [Zhang and Sabuncu, 2018].  
 116 In contrary, when  $\eta_{kj} = \eta$ , it is called symmetric noise. Under the symmetric noise assumption,  
 117 Ghosh et al. [2015] studied the functional form of loss function and concluded that by using the  
 118 symmetric loss function, one can get a global optima such that the learned model is noise tolerant.  
 119 For example, the MAE loss function is a symmetric function while the cross entropy loss function  
 120 is not. However, using MAE loss function has poor accuracy performance on classification tasks  
 121 compared with the cross entropy loss function [Zhang and Sabuncu, 2018]. One can convince  
 122 oneself with Eqn. (5) in [Zhang and Sabuncu, 2018], i.e., the cross entropy loss function enables  
 123 the optimization process weighting the sample importance while the MAE loss function considers  
 124 samples equally. Furthermore, Zhang and Sabuncu [2018] combine MAE and cross entropy loss  
 125 functions with L'Hôpital's rule, i.e.,

$$\mathcal{L}_q(f(x), j) = \frac{(1 - f_j(x)^q)}{q}, \quad (1)$$

126 where  $f(x)$  is the model,  $j$  indexes the class, and  $f_j(x)$  is the softmax output of  $j$ . Interestingly,  
 127 when  $q = 1$ ,  $\mathcal{L}_q(f(x), j)$  is a MAE loss function; while  $\lim_{q \rightarrow 0} \mathcal{L}_q(f(x), j)$  is a cross entropy loss.  
 128 Therefore, one can manipulate trade off by selecting a good hyper-parameter  $q$ . Furthermore, it  
 129 also introduces a better loss function, the truncated  $\mathcal{L}_q(f(x), j)$ , which is essentially a practically  
 130 improved version of  $\mathcal{L}_q(f(x), j)$ . However, in theory the proposed method is based on the symmetric  
 131 noise assumption [Zhang and Sabuncu, 2018], which can be quite easy to be violated. This is a  
 132 trade-off between using a stricter assumption and estimating noisy labelling mechanisms [Patrini  
 133 et al., 2017] (which is a challenge).

134 **SSL for noisy labeling and a potential solution for asymmetric noise.** Li et al. [2020] consider  
 135 the noisy label problem as a semi-supervised learning problem by finding the similarity of unlabeled  
 136 samples in semi-supervised learning and noisy labels. Suppose that we can successfully separate the  
 137 noisy and clean samples, we can treat the noisy ones as unlabeled data in semi-supervised learning,  
 138 and then leverage the success of semi-supervised learning to tackle the noisy labeling problem.

139 Firstly, by observing that the loss of clean samples tends to be lower than the noisy ones [Arazo et al.,  
 140 2019], Li et al. [2020] fit a Gaussian Mixture Model for the two components, the noisy group and  
 141 the clean one. Then given a loss, it can be inferred whether the sample is a noisy one or a clean  
 142 one. Consequently, following the mentioned idea, semi-supervised learning methods are applied  
 143 to such a separated dataset. Moreover, Li et al. [2020] consider the influence of asymmetric noise  
 144 in the supervised learning phase. Because the bias introduced by the asymmetric noise can lead  
 145 to severe consequences (confirmation errors). [Li et al., 2020] added a negative entropy penalty  
 146 term  $-\mathcal{H} = \sum_j f_j(x) \log f_j(x)$  for an input  $x$  in the cross-entropy loss function at the beginning  
 147 of training to avoid over-confident prediction, which works well empirically. To further reduce  
 148 the influence of the confirmation error introduced by the symmetric noise, it uses the MixMatch  
 149 [Berthelot et al., 2019] procedure to train two independent DNNs and attractively exchange datasets  
 150 with each other for filtering errors made by the other one. This is actually an ensemble method, which  
 151 reduces the random noise in the prediction, especially in the presence of symmetric labelling noise.

152 **Model bias in SSL.** Kurakin et al. [2020] propose a distribution alignment method utilizing a  
 153 principle introduced by Bridle et al. [1992]. It formulates an ideal classifier which maximizes the  
 154 mutual information of model inputs and model outputs. Furthermore, it argues that the second term  
 155 of mutual information encourages a model to output with low entropy and high confidence, while  
 156 another one encourages equal frequency across the entire training set as shown in

$$\begin{aligned} \mathcal{I}(y; x) &= \iint \log \frac{p(y, x)}{p(y)p(x)} dy dx \\ &= \mathcal{H}[\mathbb{E}[p(y | x; \theta)]] - \mathbb{E}_x[\mathcal{H}[p(y | x; \theta)]], \end{aligned} \quad (2)$$

157 where  $\theta$  is the model parameters. As what Kurakin et al. [2020] said, when the marginal distribution  
 158 of a training dataset labels is not uniformly distributed, it is not proper to regularize the frequency. In  
 159 our work, to deal with such case, we augment the training dataset and make the labels of labeled data  
 160 in each batches to be uniformly distributed.

## 161 3 Methods

### 162 3.1 FixMatch

163 As one of the SSL methods, FixMatch [Sohn et al., 2020] leverages labeled data and introduces prior  
 164 knowledge about unlabeled data in the training process. For labeled data, FixMatch simply uses the  
 165 cross entropy loss function for a batch,

$$l_s = \frac{1}{B} \sum_{b=1}^B H(y_b, f(\alpha(x_b))), \quad (3)$$

166 where  $B$  is the number of labeled data in a batch,  $x_b$  is a labeled sample,  $y_b$  is the label, and  $\alpha(\cdot)$  is  
 167 weak augmentation. However, due to limited number of labeled samples, the performance of such  
 168 DNN is not ideal. Therefore, the question is how to make a good use of the sufficient unlabeled data  
 169 to improve the performance? Ideally, the performance can be close to the DNN trained with the fully  
 170 labeled dataset.

171 FixMatch considers the consistency of model prediction on the unlabeled data with weak and strong  
 172 augmentation (the augmentation methods are introduced in Sec. 4). It first uses the model to predict  
 173 pseudo labels for unlabeled data and then compute the loss of unlabeled data with the pseudo labels  
 174 and the consistency regularization. The loss function for the unlabeled samples  $u_b$  is

$$l_u = \frac{1}{\mu B} \sum_{b=1}^{\mu B} \mathbf{1}(\max_y (f(\alpha(u_b))) \geq \tau) H(\hat{y}_b, f(\mathcal{A}(u_b))), \quad (4)$$

175 where  $\mu B$  is the number of unlabeled data in a batch,  $\hat{y}_b := \arg \max_y p(y | \alpha(u_b); \theta_f)$  is the pseudo  
 176 label of  $u_b$ ,  $\theta_f$  is the neural network parameters of function  $f$ , and  $\mathcal{A}(\cdot)$  is the strong augmentation.  
 177 Note that to make pseudo labels reliable to be used, FixMatch considers the pseudo labels in the loss  
 178 function only if the prediction has a higher probability than  $\tau$ . Next, together with the cross entropy  
 179 loss of labeled data, the loss function of FixMatch is  $l_s + \lambda_u l_u$ .

180 **3.2 Investigation of noisy (pseudo) labels and confirmation errors of FixMatch**

181 **Noisy pseudo labels and confirmation errors in FixMatch.** A main issue of "Match"-based SSL  
 182 methods is confirmation errors. Since FixMatch is trained on batches with both labeled and unlabeled  
 183 data, it is very likely to make prediction errors at the beginning of the training. When the model  
 184 makes wrong predictions of labeled data, since we have their ground-truth labels, the model can  
 185 become better with the loss for labeled data. But when it comes to unlabeled samples, since we don't  
 186 have the ground-truth labels, the model uses the confident pseudo labels as the labels for training.  
 187 In this case, if the pseudo-labels are noisy, the model can fit such errors and become biased. In  
 188 the next batch, it can generate more wrong pseudo-labels with higher confidence. Moreover, the  
 189 consistency regularization can keep reinforcing the model to fit such wrongly labeled data. Finally, it  
 190 demonstrates a biased model with a poor performance on generalization and robustness. Therefore,  
 191 noise in the pseudo labels can lead to confirmation errors in FixMatch.

192 Both asymmetric noise and symmetric noise in pseudo labels can lead to confirmation errors, but in  
 193 general asymmetric noise is more harmful and harder to deal with. For example, to reduce the impact  
 194 of symmetric noise and get an unbiased model, one can use ensembling methods like [Li et al., 2020]  
 195 to train multiple DNNs at the same time; however, this can fail in the presence of asymmetric noise.  
 196 In this work, we focus on asymmetric noise and one can simply extend the method to deal with the  
 197 influence of symmetric noise with ensemble methods.

198 **Investigation with class-balanced batches.** To check whether there exist confirmation errors, we  
 199 need to check that during the training process errors are reinforced by the model. Moreover, to  
 200 see the asymmetric noise in the pseudo labels, we need to check that in the training phase whether  
 201 FixMatch predicts a certain class of unlabeled data into certain other classes. Thus, these require us  
 202 to investigate the performance of FixMatch at each batch and check the pseudo labelling performance  
 203 regarding asymmetric noise in the pseudo labels. However, in [Sohn et al., 2020], a batch is not  
 204 necessary to contain all the classes of training dataset and it can contain different classes with different  
 205 numbers. Therefore, the performance of pseudo labelling regarding asymmetric noise inherits the  
 206 randomness of batch composition, which makes the investigation conclusion unreliable.

207 To deal with this issue, we reconstructed the batch structure which requires each batch to contain an  
 208 equal number of images for all the classes on both labeled and unlabeled data, called Balanced-Class  
 209 (BC) batches. With such batches, we can fairly check the performance of pseudo labelling in each  
 210 batch how many errors are made when the model predicts each class and whether it tends to label  
 211 a class as other certain classes causing asymmetric noise. Note that without further introducing  
 212 regularization, BC batches on their own cannot improve the performance of FixMatch, which has  
 213 indistinguishable results without BC as shown in Sec. 5.3.

214 Furthermore, we leverage the reconstructed batch structure to regularize the training process for  
 215 reducing the noise in pseudo labels and improving the performance. With the reconstructed batches,  
 216 we know that the class of labeled data<sup>1</sup> is uniformly distributed, thus we can regularize the output  
 217 of labeled data with the negative entropy loss of the prediction frequency. In this way we force the  
 218 output of labeled data to be uniformly distributed. Potentially this can regularize the asymmetric noise  
 219 in the labeled data because the output class distribution is not likely to be uniformly distributed in the  
 220 presence of asymmetric noise. Consequently, it can reduce the asymmetric noise in pseudo labels  
 221 because the prediction on both labeled and unlabeled data uses the same network which is unlikely to  
 222 have different prediction behavior. Therefore, we add an equal-frequency entropy regularization to  
 223 the loss function, which is

$$\begin{aligned}
 l' &= l'_s + \lambda_u l_u, \\
 l'_s &= l_s - \lambda_{ef} \mathcal{H}(\mathbb{E}_{x_b}[f(\alpha(x_b))]) \\
 &= l_s + \lambda_{ef} \sum_{j=1}^c \left\{ \left( \frac{1}{B} \sum_{b=1}^B f_j(\alpha(x_b)) \right) \log \left( \frac{1}{B} \sum_{b=1}^B f_j(\alpha(x_b)) \right) \right\},
 \end{aligned} \tag{5}$$

<sup>1</sup>In fact, the class of both labeled and unlabeled data are equally distributed in reconstructed batches, but it is unrealistic to use the prior knowledge about labels of unlabeled data. Although it is fine for "debugging" the training behavior of FixMatch, when aiming at improving the performance of FixMatch, we cannot use the information about labels of unlabeled data, because it is very likely to have unbalanced classes of unlabeled data in practice. Then it makes no sense to regularize the outputs of unlabeled data in the training phase.

224 where  $c$  is the number of classes and  $\lambda_{ef}$  is a hyperparameter. We also consider the confidence  
 225 entropy loss regularization which can avoid over-confident prediction,

$$\begin{aligned}
 l''_s &= l_s - \lambda_{ce} \mathbb{E}_{x_b} [\mathcal{H}(f(\alpha(x_b)))] \\
 &= l_s + \lambda_{ce} \frac{1}{B} \sum_{b=1}^B \left\{ \sum_{j=1}^c f_j(\alpha(x_b)) \log(f_j(\alpha(x_b))) \right\}, \\
 l'' &= l''_s + \lambda_u l_u.
 \end{aligned} \tag{6}$$

Note that since the loss function (6) aims for avoiding over-confident predictions, it seems to be fine to regularize the unlabeled data as well. However, we cannot do that for the same reason as the loss function (5) which has been discussed in the footnote. Because  $-\mathcal{H}(\cdot)$  is a convex function, we have the Jensen's inequality

$$-\mathcal{H}(\mathbb{E}_{x_b}[f(\alpha(x_b))]) \leq -\mathbb{E}_{x_b}[\mathcal{H}(f(\alpha(x_b)))].$$

226 In other words, confidence entropy regularization can implicitly regularize the equal frequency of the  
 227 data labels. Therefore, with the same reason, we should only apply it to the labeled data of which  
 228 label distribution is under our control with augmentation.

## 229 4 Data Preprocessing and Augmentation

230 FixMatch requires a weak augmentation  $\alpha(\cdot)$  and a strong augmentation  $\mathcal{A}(\cdot)$ . For the weak augmen-  
 231 tation, we randomly flip an image with probability 0.5 as [Sohn et al., 2020] and translate an image  
 232 up to 12.5% with probability 0.5<sup>2</sup>. For the strong augmentation, FixMatch uses either RandAugment  
 233 (RA) [Cubuk et al., 2020] or CTAugment [Kurakin et al., 2020] for their experiments. However, we  
 234 use RA for our experiments with the maximum magnitude 10 (same as the official experiment setup)  
 235 and 2 randomly selected operations per image.

236 Due to the limitation of computational resources, we examine the reproducibility of [Sohn et al.,  
 237 2020] on the dataset CIFAR-10 [Krizhevsky et al., 2009]. In CIFAR-10, there are 50000 training  
 238 data and 10000 testing data. We take 5000 training data as the validation dataset. Then we use the  
 239 remaining training dataset to make labeled and unlabeled datasets and augment both datasets into the  
 240 same target number as in [Sohn et al., 2020]. After augmentation, we have  $2^{13}$  labeled images and  
 241  $2^{13} \times 7$  unlabeled images for the CIFAR-10 training dataset.

## 242 5 Experiment

243 In the reproducibility experiments, we re-implement FixMatch from scratch using PyTorch and  
 244 reproduce the essential experiments in the original paper with the similar results. We use the  
 245 hyperparameters ( $\lambda_u = 1$ ,  $\eta = 0.03$ ,  $\beta = 0.9$ ,  $\tau = 0.95$ ,  $\mu = 7$ ,  $B = 64$ ,  $K = 2^{20}$ ) given by  
 246 [Sohn et al., 2020] and focus on reproducing the performance on CIFAR-10 (Sec. 4.1 of [Sohn et al.,  
 247 2020]) and barely supervised learning experiments (Sec. 4.4 of [Sohn et al., 2020]). Besides the early  
 248 introduced hyper-parameters, we use SGD with  $\beta = 0.9$  for training the model, and the learning  
 249 rate is decay with  $\eta \cos(\frac{7\pi k}{16K})$ , where  $K$  is the total time step and  $k$  is the current time step. Each  
 250 experiment takes around 68 hours on a single V100. And all the error rates is generated from EMA  
 251 (exponential moving average) of models in the SGD training trajectory.

252 Then, we investigate confirmation errors of "Match"-based SSL methods to see whether there exists  
 253 such error and asymmetric noise of pseudo labels in FixMatch with the official experiment setup, i.e.  
 254 unbalanced batches, in [Sohn et al., 2020]. Next, we examine the existence of confirmation errors  
 255 and asymmetric noise for FixMatch again in a more reliable way using re-constructed batches as  
 256 introduced in Sec. 3. Furthermore, we respectively add the equal-frequency entropy regularization  
 257 and confidence entropy regularization on the labeled training data in the loss function and compare  
 258 with the baseline FixMatch without entropy regularization on the BC batches. Finally, we add  
 259 asymmetric noise to the labeled data in the training dataset and compare the performance of baseline  
 260 FixMatch and FixMatch with different entropy regularization.

<sup>2</sup>Here, [Sohn et al., 2020] didn't indicate what probability they use for the translation.

261 **5.1 Reproducibility**

262 **CIFAR-10.** We reproduced the experiments on CIFAR-10 with 40, 250, 4000 labeled data and  
263 5000 validation samples as the official implementation of FixMatch<sup>3</sup>. But due to the limitation of  
264 computational resources, we didn't reproduce 5 "folds". Thus, our result based on 1 fold doesn't have  
265 the standard deviation. Our model uses the Wide ResNet-28-2 [Zagoruyko and Komodakis, 2016]  
266 with leaky ReLU activation function. Our results are shown in Table 1 which is comparable to the  
267 performance in [Sohn et al., 2020].

Table 1: Error rates for CIFAR-10 on test data. FixMatch (RA) uses RandAugment [Cubuk et al., 2020]. BC means that the experiment uses balanced-class batches as introduced in Sec. 3. We use the experiment with BC and RA as a comparison baseline results for the investigation in Sec. 5.3.

Method	CIFAR-10		
	40 labels	250 labels	4000 labels
Official FixMatch (RA)	13.81 ± 3.37	5.07 ± 0.65	4.26 ± 0.05
Ours (RA)	10.04	5.29	4.36

268 **Barely supervised learning.** We also reproduce the one example per class experiment. [Sohn  
269 et al., 2020] hypothesize that the repressiveness of the chosen labeled data influences the results  
270 significantly. Since there are only one/few samples per class, this hypothesis is reasonable intuitively.  
271 Then, Sohn et al. [2020] categorized the training dataset into eight levels of "prototypicality", i.e.,  
272 representative of the underlying class and then ordered the training samples by their "prototypicality".  
273 With the same hyperparameters, the model is trained with 10 provided most representative labeled  
274 data under Random Augment. The accuracy is 84.41% compared with the official implementation: a  
275 median of 78% accuracy and a maximum of 84% accuracy.

276 **5.2 Ablation studies**

277 The ablation studies are based on FixMatch with 250 labels using CTAugment.

278 **Study for Confidence threshold.** We performance the ablation studies for confidence threshold.  
279 Due to the limited computation resource, we hypothesize that experiments with lower confidence  
280 threshold will achieve worse performance and explore more values around the optimal value of  
281 confidence threshold, 0.95 chosen by the authors. Thus our examined threshold value is between  
282 0.7 to 0.98. As shown in Figure1 (c), the error rate is between 6.54% and 6.19% and the highest  
283 performance is under the threshold 0.98.

284 **Ratio of unlabeled data.** We perform FixMatch under different ratios of unlabeled data. Figure1  
285 (d) shows the error rate which is decreasing when the ratio of unlabeled data is higher. A significant  
286 increase of the accuracy happens using a large number of unlabeled data. The results show the  
287 consistency with the finding in the original paper.

288 **5.3 Investigation on confirmation errors and asymmetric noisy (pseudo) labels**

289 In this section, we show the investigation on confirmation errors and asymmetric noise in labels and  
290 pseudo labels and whether the entropy regularization in loss functions (5) and (6) can deal with them  
291 and improve the performance of FixMatch. The training dataset contains 150 labeled data before  
292 augmentation and each BC batch in the training phase contains images with uniformly distributed  
293 classes.

294 **Existence of asymmetric noise and confirmation errors in pseudo labels.** We examine the  
295 existence of asymmetric noise in pseudo labels by checking the confusion matrix of the prediction  
296 of unlabeled data in different batches. Top figures of Figure 2 show the confusion matrices in the  
297 experiments without using BC batches. We find that asymmetric noise appears in a random manner,  
298 which is as our expectation as analyzed in Sec. 3. The stochastic behavior is inherited from the

<sup>3</sup>The official implementation: <https://github.com/google-research/fixmatch>. From the reproducibility and readability, the official code is not a valid submission.

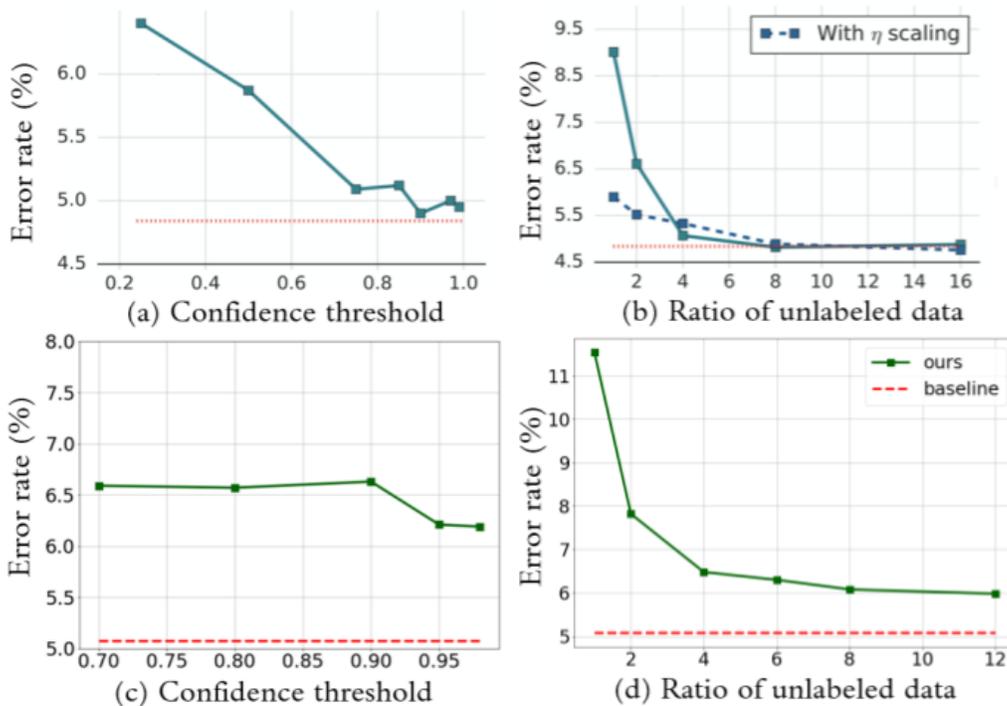


Figure 1: Plots of various ablation studies on FixMatch compared with those reported in the paper. (a) Varying the confidence threshold for pseudo-labels in the original paper. (b) Varying the ratio of unlabeled data ( $\mu$ ) in the original paper. (c) Varying the confidence threshold for pseudo-labels based on our implementation. (d) Varying the ratio of unlabeled data ( $\mu$ ) based on our implementation.

299 randomness of batch composition. Next, we evaluate the asymmetric noise with BC batches, which  
 300 is a more reliable way as mentioned in Sec. 3. We found that there exists consistent asymmetric  
 301 noise, which leads to the confirmation errors, i.e., the model always tends to wrongly predict certain  
 302 images into certain classes as shown in bottom figures of Figure 2. Moreover, the accuracy of our  
 303 implementation is 93.6% without BC batches and 93.8% with BC batches, which shows that using  
 304 BC batches has rarely influence on the model performance compared with the one without BC  
 305 batches.

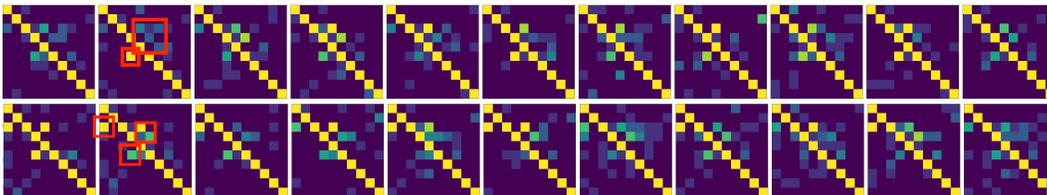


Figure 2: Confusion matrices of the confident prediction on unlabeled data with different batch structures. Confusion matrices are plotted every 100 training steps in the 1st epoch (1024 steps). The **top** matrices are from the experiments without BC, and the **bottom** matrices are the ones with BC. The red areas represent the asymmetric noise in the pseudo labels. The bottom matrices have a stable and smooth transition while the top matrices have a fluctuating transition in the red areas. The yellow color represents larger value and the darker green color represents smaller values.

306 **Equal-frequency and confidence entropy regularization on the labeled data.** Due to limitation  
 307 of the computational resources, we didn't explicitly run grid search for the hyperparameters in the  
 308 Equal-Frequency (EF) loss function (5) and Confidence-Entropy (CE) loss function (6). Instead,  
 309 we found that for the baseline method the training loss is around 0.2. We then compute the equal-

310 frequency entropy loss for the ideal scenario, equal frequency for all classes, which is  $0.1 \times \ln 0.1 \approx 2$ .  
 311 We decide to try the hyperparameter  $\lambda_{ce}$ ,  $\lambda_{ef} \leq 0.1$  to avoid making the entropy regularization loss  
 312 dominate the loss value. Then, we do a hyper-parameter search for the loss function (5) and (6). For  
 313 all experiments in this experiment, we used cosine function decay for the parameters  $\lambda_{ce}$  and  $\lambda_{ef}$ ,  
 314 which starts with value 1 and ends with value 0 in the training phase. We find that using the loss  
 315 function (6) can achieve a better accuracy performance 94.01%. Moreover, as an advantage, using the  
 316 confidence entropy regularization can reduce the asymmetric noise as shown in the bottom confusion  
 317 matrices of Figure 3. As for the equal-frequency entropy regularization, it has a better accuracy,  
 318 93.85%. Moreover, the equal-frequency entropy regularization can penalize the asymmetric noise,  
 319 which may transform it into symmetric noise as shown in the middle confusion matrices of Figure 3.  
 320 Note that there are plenty of ways to deal with symmetric noise, which is much easier to handle.

Table 2: Error rates on testing data using the loss function (5) and (6). The experiments use 150 labeled data and CTA for training. The first column is the results without BC batch and the second column is the baseline result without using EF or CE regularization.

Entropy regularization	noBC+Null	BC+Null	BC+CE	BC+EF
$\lambda_{ce}/\lambda_{ef}$	0	0	0.1	0.1
Error rate	6.4	6.2	<b>5.99</b>	6.15

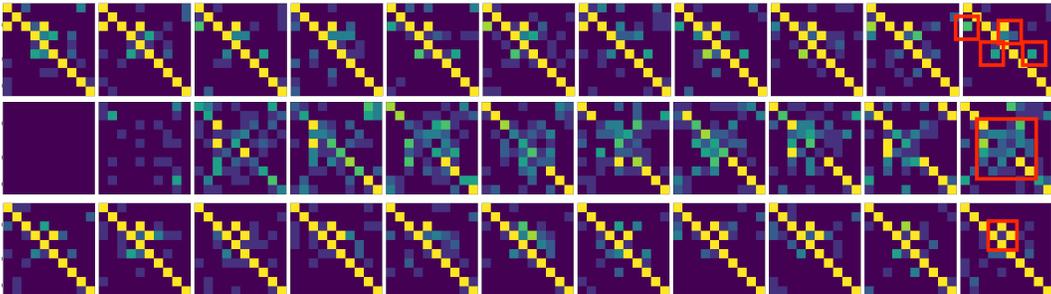


Figure 3: Confusion matrices of the confident prediction on unlabeled data with BC batches using loss functions (4) without entropy regularization at **top**, (5) with equal-frequency entropy regularization in the **middle**, and (6) with confidence entropy regularization at **bottom**. Confusion matrices are plotted every 100 training steps in the 1st epoch (1024 steps). The red areas represent the asymmetric/symmetric noise in the pseudo labels. The yellow color represents larger value and the darker green color represents smaller values.

321 **Equal-frequency and confidence entropy regularization on the labeled data containing asym-**  
 322 **metric noise.** In this experiment, we use RA data augmentation and manually add asymmetric  
 323 noise to the labeled data in the training dataset to compare how FixMatch with different loss functions  
 324 performs in the presence of asymmetric noise in the labeled data. We respectively select 3 images  
 325 from class 0 and class 1 in the validation dataset. Then, for the labeled data in the training dataset,  
 326 we keep the labels unchanged and replace 3 images in class 2 with the 3 images in class 0. Similarly  
 327 we replace 3 images in class 3 with the 3 images in class 1. In this way, the only difference with  
 328 the previous experiments in this section is that our final validation dataset has 4994 images and the  
 329 labeled data in the training dataset contain asymmetric noise. Table 3 shows error rates on 6 runs with  
 330 different random seeds. In the presence of asymmetric noise in labeled training data, all proposed  
 331 methods perform better than the baseline method, in which FixMatch with BC batches decreases the  
 332 average error rate from 8.6 to 7.37, and the combination of confidence-entropy regularization and BC  
 333 batches further lowers the error rate to 6.98.

## 334 6 Challenges

335 It is not clear how many steps are there in each epoch. First the paper only states the total steps  
 336  $K = 2^{20}$  and the composition of one batch ( $B$  labeled samples and  $\mu B$  unlabeled samples). And the  
 337 official code indicates there are  $2^{16}$  labeled images observed by the model per epoch and a total of

Table 3: Error rates of FixMatch methods in the presence of asymmetric noise in labeled training data augmented by RA: The baseline method ( $\lambda = 0$ ); The method ( $\lambda = 0$ ) with BC batches; the method with confidence-entropy regularization ( $\lambda_{ce} = 0.1$ ) and BC batches; the method with equal-frequency regularization ( $\lambda_{ef} = 0.1$ ) and BC batches.

	$\lambda = 0(\text{noBC})$	$\lambda = 0(\text{BC})$	$\lambda_{ef} = 0.1(\text{BC})$	$\lambda_{ce} = 0.1(\text{BC})$
Error rates on test data	$8.6 \pm 2.81$	$7.37 \pm 2.05$	$7.95 \pm 2.2$	<b><math>6.98 \pm 1.83</math></b>

338  $2^{26}$  images observed which suggests that there are  $2^{12}$  updates per epoch and  $2^{19}$  updates in total.  
 339 And this is not consistent with the total update steps  $K$  stated in the paper. When performing weak  
 340 augmentations to the input data, the probability for randomly translating images is not specified. And  
 341 it also remains unclear the ‘5 different folds’ mentioned in the paper, we are guessing it is a kind of  
 342 cross validation while there is not too much evidence supporting this neither in the paper nor in the  
 343 official code.

344 The paper doesn’t contain sufficient details to reproduce all the experiments. Thus, it is necessary to  
 345 look for details about reproducing the experiments in the official code. We have not optimized or  
 346 tuned the hyperparameters, and all the hyperparameters are the same as those mentioned in the paper.  
 347 Compared to the average error rates in the original paper, the reproduced results have a reasonable  
 348 good performance on a larger number of labeled data (4000/250 labels) and better but also reasonable  
 349 performance on fewer labeled data (40/10 labels) since the variance of error rates over 5 different  
 350 folds for CIFAR-10 with 40 labels is 3.35%. Moreover, to compare with the results of ablation studies  
 351 in the original paper, we also implement CTAugment, which supports a learnable magnitude. While  
 352 we failed to confirm the result that CTAugment behaves better than RandAugment on CIFAR-10. We  
 353 hypothetically guess this is because it could affect the consistency regularization because of different  
 354 levels of distortions controlled by magnitude.

## 355 7 Conclusion

356 In this work, we study and reimplement FixMatch from scratch. We reproduced essential experiments,  
 357 included the model performance on CIFAR 10, barely supervised learning, and ablation studies.  
 358 Experimental results show that our implementation achieves similar performance as the original  
 359 FixMatch results, which supports that FixMatch outperforms semi-supervised learning benchmarks  
 360 and that the author’s choices with respect to those ablations were experimentally sound. We also  
 361 confirmed the existence of confirmation errors in pseudo labels by checking the prediction confusion  
 362 matrix of unlabeled data in different training stages. We adapted the training batch structure to be  
 363 composed of equal number of images in each class, which enable us to stably and reliably check the  
 364 the asymmetric noise in the training phase. Based on the reconstructed batch structure, we used the  
 365 equal-frequency and confidence entropy regularization in the loss function, and theoretically show  
 366 their relation. The experiments indicate that these entropy regularization can reduce the asymmetric  
 367 noise in pseudo labels and improves the performance of FixMatch in the presence of training labels  
 368 with asymmetric noise.

## 369 8 Ethical consideration

370 The bias in the collected dataset is a serious problem when applying machine learning methods to  
 371 the real-world scenarios. For example, applying machine learning methods to making automated  
 372 decision-making systems for criminal prediction, university admission or recruitment. In these cases,  
 373 we may very likely collect a dataset containing certain bias due to the historical reason or selection  
 374 bias in the data collection process. If a model cannot deal with such bias in the dataset, it may inherit  
 375 in the model by focusing on the unrelated or wrong relations in the dataset. Consequently, the model  
 376 can make biased decision which can disadvantage a certain group of people and may even diminish  
 377 this group in the society.

378 Unfortunately, FixMatch cannot only be influenced by the noise in the label of a training dataset, but  
 379 also it can make confirmation errors causing a biased model even when the dataset itself is unbiased.  
 380 To deal with such issue, this work focuses on the asymmetric noise in the data labels and pseudo

381 labels, which can lead to severe confirmation error and the biased model. And then, we applied  
382 different methods to reduce such noise in pseudo labels and reduce its impact on the model.

### 383 **References**

- 384 E. Arazo, D. Ortego, P. Albert, N. E. O’Connor, and K. McGuinness. Unsupervised label noise  
385 modeling and loss correction. *arXiv preprint arXiv:1904.11238*, 2019.
- 386 S. Ben-David, T. Lu, and D. Pál. Does unlabeled data provably help? worst-case analysis of the  
387 sample complexity of semi-supervised learning.
- 388 D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel. Mixmatch: A  
389 holistic approach to semi-supervised learning. In *Advances in Neural Information Processing*  
390 *Systems*, pages 5049–5059, 2019.
- 391 J. S. Bridle, A. J. Heading, and D. J. MacKay. Unsupervised classifiers, mutual information  
392 and ‘phantom targets. In *Advances in neural information processing systems*, pages 1096–1101,  
393 1992.
- 394 O. Chapelle, B. Schölkopf, and A. Zien. Semi-supervised learning. adaptive computation and machine  
395 learning. *MIT Press, Cambridge, MA, USA. Cited in page (s), 21(1):2*, 2010.
- 396 E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le. Randaugment: Practical automated data augmentation  
397 with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
398 *and Pattern Recognition Workshops*, pages 702–703, 2020.
- 399 Z. Ghahramani. Keynote: Machine learning and a.i. at uber. [https://www.youtube.com/watch?](https://www.youtube.com/watch?v=4XTv5qgugCk&feature=youtu.be)  
400 [v=4XTv5qgugCk&feature=youtu.be](https://www.youtube.com/watch?v=4XTv5qgugCk&feature=youtu.be), 2020.
- 401 A. Ghosh, N. Manwani, and P. Sastry. Making risk minimization tolerant to label noise. *Neurocom-*  
402 *puting*, 160:93–107, 2015.
- 403 A. Krizhevsky et al. Learning multiple layers of features from tiny images. 2009.
- 404 A. Kurakin, C. Raffel, D. Berthelot, E. D. Cubuk, H. Zhang, K. Sohn, and N. Carlini. Remixmatch:  
405 Semi-supervised learning with distribution matching and augmentation anchoring. 2020.
- 406 A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov,  
407 M. Mallocci, T. Duerig, et al. The open images dataset v4: Unified image classification, object  
408 detection, and visual relationship detection at scale. *arXiv preprint arXiv:1811.00982*, 2018.
- 409 J. Li, R. Socher, and S. C. Hoi. Dividemix: Learning with noisy labels as semi-supervised learning.  
410 *ICLR*, 2020.
- 411 G. Patrini, A. Rozza, A. Krishna Menon, R. Nock, and L. Qu. Making deep neural networks robust  
412 to label noise: A loss correction approach. In *Proceedings of the IEEE Conference on Computer*  
413 *Vision and Pattern Recognition*, pages 1944–1952, 2017.
- 414 B. Schölkopf, D. Janzing, J. Peters, E. Sgouritsa, K. Zhang, and J. Mooij. On causal and anticausal  
415 learning. *arXiv preprint arXiv:1206.6471*, 2012.
- 416 K. Sohn, D. Berthelot, C.-L. Li, Z. Zhang, N. Carlini, E. D. Cubuk, A. Kurakin, H. Zhang, and  
417 C. Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *arXiv*  
418 *preprint arXiv:2001.07685*, 2020.
- 419 N. Tajbakhsh, J. Y. Shin, S. R. Gurudu, R. T. Hurst, C. B. Kendall, M. B. Gotway, and J. Liang.  
420 Convolutional neural networks for medical image analysis: Full training or fine tuning? *IEEE*  
421 *transactions on medical imaging*, 35(5):1299–1312, 2016.
- 422 A. Tarvainen and H. Valpola. Weight-averaged consistency targets improve semi-supervised deep  
423 learning results.

- 424 J. von Kügelgen, A. Mey, M. Loog, and B. Schölkopf. Semi-supervised learning, causality, and  
425 the conditional cluster assumption. volume 124 of *Proceedings of Machine Learning Research*,  
426 pages 1–10, Virtual, 03–06 Aug 2020. PMLR. URL [http://proceedings.mlr.press/v124/  
427 kugelgen20a.html](http://proceedings.mlr.press/v124/kugelgen20a.html).
- 428 S. Zagoruyko and N. Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- 429 Z. Zhang and M. Sabuncu. Generalized cross entropy loss for training deep neural networks with  
430 noisy labels. In *Advances in neural information processing systems*, pages 8778–8788, 2018.