

MarginMatch: Improving Semi-Supervised Learning with Pseudo-Margins

Tiberiu Sosea

Cornelia Caragea

University of Illinois Chicago

tsosea2@uic.edu

cornelia@uic.edu

What's semi-supervised learning?



Few labeled data



Numerous unlabeled data

Development of Different SSL Methods



$$\ell_u = \frac{1}{\mu B} \sum_{b=1}^{\mu B} \mathbb{1}(\max(q_b) \geq \tau) H(\hat{q}_b, p_m(y | \mathcal{A}(u_b)))$$

Motivation

How to make unlabeled sample selection more adaptive?



									
FixMatch									
Predicted:	onion	elephant	fossa	green pepper	pop art	crowd	firefighter	horse	crowd
Actual:	bell pepper	camel	cougar	handrail	poncho	uniform	volleyball	bison	meat market
									
FlexMatch									
Predicted:	screen	pyramid	decoration	scale	computer	carpet	cabbage	tower	screen
Actual:	stopwatch	obelisk	socks	parking meter	heater	teddy bear	cauliflower	torch	ipod

Figure 1. Incorrect pseudo-labels propagated until the end of the training process for FixMatch and FlexMatch on ImageNet.

Confirmation Bias

If the initial predictions are **biased** or **contain errors**, the model would tend to be **overconfident** in its incorrect predictions.

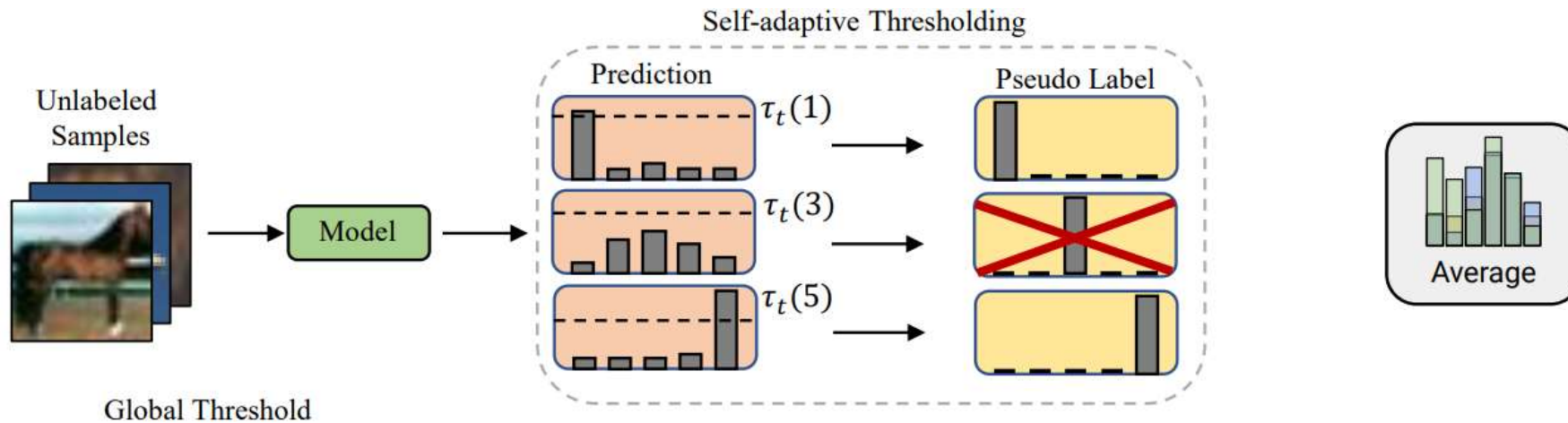


Overconfident

Even when the **high-confidence threshold** is used in FixMatch can result in wrong pseudo-labels.

Motivation

One Strike and you're out?



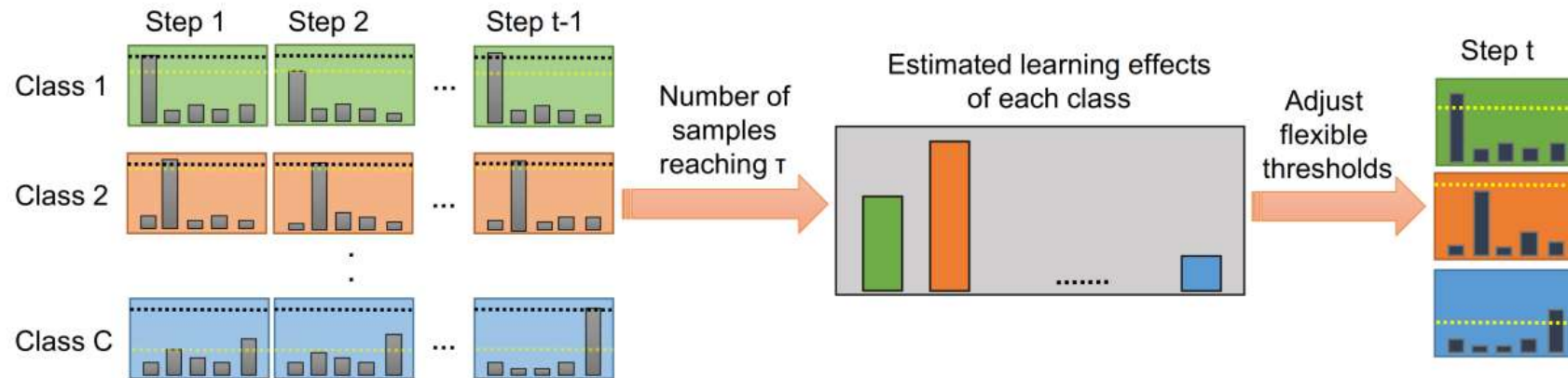
Instead of using only the model's confidence on an unlabeled example at **an arbitrary iteration** to decide if the example should be masked or not, MarginMatch **analyzes the behavior** of the model on the pseudo-labeled examples **as the training progresses**, to ensure **low quality predictions** are masked out.

Instead of using *One Strike and you're out* rule

How to measure prediction quality more accurately?

Prerequisite

FlexMatch



$$\sigma_t(c) = \sum_{n=1}^N \mathbb{1}(\max(p_{m,t}(y|u_n)) > \tau) \cdot \mathbb{1}(\arg \max(p_{m,t}(y|u_n)) = c).$$

Reflects the learning effect of class c at time step t

$$\mathcal{L}_{u,t} = \frac{1}{\mu B} \sum_{b=1}^{\mu B} \mathbb{1}(\max(q_b) > \mathcal{T}_t(\arg \max(q_b))) H(\hat{q}_b, p_m(y|\Omega(u_b))),$$

$$\beta_t(c) = \frac{\sigma_t(c)}{\max_c \sigma_t},$$

$$\mathcal{T}_t(c) = \beta_t(c) \cdot \tau.$$

Instead of using only the model's **current belief**, MarginMatch monitors the training dynamics of unlabeled data over the iterations by investigating the margins:

$$\text{PM}_c^t(\hat{x}) = z_c - \max_{c' \neq i} (z_{c'})$$

We average all the margins with respect to c from the first iteration until t and obtain the average pseudo-margin (APM) as follows:

$$\text{APM}_c^t(\hat{x}) = \frac{1}{t} \sum_{j=1}^t \text{PM}_c^j(\hat{x})$$

If over the iterations, the model predictions do **not agree frequently** with the pseudo-label c from iteration t on the predicted label, the APM for class c will have a **low** value.

Old pseudo-margins eventually become deprecated due to the **large** number of **unlabeled data**, so we use EMA to place **more importance** on **recent iterations**.

$$\text{APM}_c^t(\hat{x}) = \text{PM}_c^t(\hat{x}) * \frac{\delta}{1+t} + \text{APM}_c^{t-1}(\hat{x}) * \left(1 - \frac{\delta}{1+t}\right)$$

Illustration of Pseudo-margin (PM) Metric

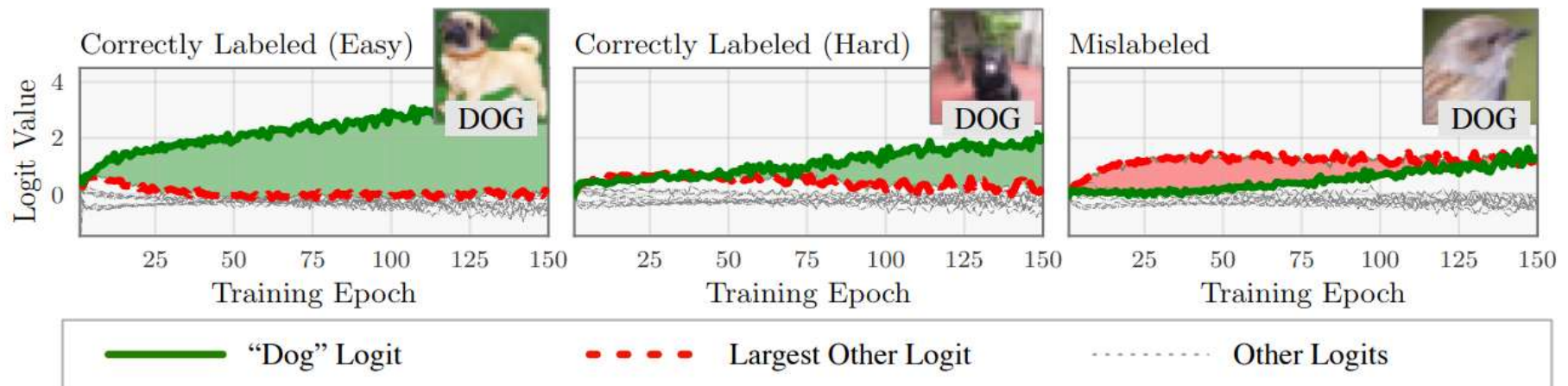


Figure 2: Illustration of the *Area Under the Margin* (AUM) metric. The graphs display logit trajectories for easy-to-learn dogs (left), hard-to-learn dogs (middle), and BIRDS mislabeled as DOGS (right). (Each plot's logits are averaged from 50 CIFAR10 training samples, 40% label noise.) AUM is the shaded region between the DOG logit and the largest other logit. Green/red regions represent positive/negative AUM. Correctly-labeled samples have larger AUMs than mislabeled samples.

Identifying Mislabeled Data using the Area Under the Margin Ranking

The margin at epoch t captures how much larger the (potentially incorrect) assigned logit is than all other logits:

$$M^{(t)}(\mathbf{x}, y) = \overbrace{z_y^{(t)}(\mathbf{x})}^{\text{assigned logit}} - \overbrace{\max_{i \neq y} z_i^{(t)}(\mathbf{x})}^{\text{largest other logit}}.$$

Average a sample's margin measured at each training epoch, a metric we refer to as area under the margin (AUM):

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

$$\mathcal{L}_u = \sum_{i=1}^{\nu B} \mathbf{1} \left(\overbrace{\text{AM}_{\hat{p}_\theta(y|\pi(\hat{x}_i))}^t(\hat{x}_i) > \gamma^t}^{\text{APM(AUM)}} \right) \mathbf{1} \left(\overbrace{\max(p_\theta(y|\pi(\hat{x}_i))) > \mathcal{T}_{\hat{p}_\theta(y|\pi(\hat{x}_i))}^t H(\hat{p}_\theta(y|\pi(\hat{x}_i)), p_\theta(y|\Pi(\hat{x}_i)))}^{\text{FlexMatch}} \right)$$

How to Estimate APM(AUM) Threshold?

Mislabeled Data Can Be Identified by Setting a Threshold

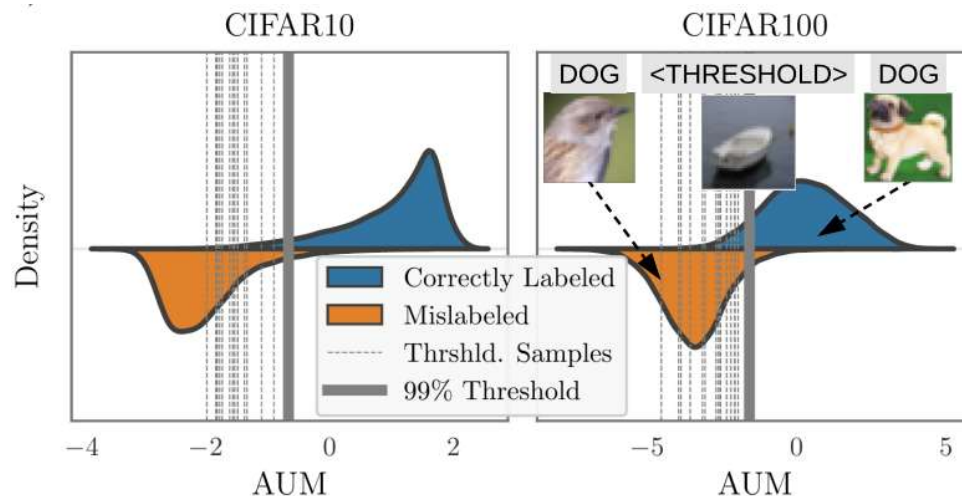


Figure 3: Illustrating the role of *threshold samples* on CIFAR10/100 with 40% mislabeled samples. Histograms of AUMs for correctly-labeled (blue) and mislabeled samples (orange). Dashed lines represent the AUM values of threshold samples. The 99th percentile of threshold AUMs (solid gray line) separates correctly- and mislabeled data.

With access to a trusted **validation set**, a threshold can be learned through a hyperparameter sweep.

But how to learn a threshold without validation data?

Randomly sample a subset of unlabeled examples from U to create **erroneous(threshold)** examples E , which are assigned to an inexistent (or virtual) class $C + 1$.

Compute APM_{C+1}^t to choose the APM threshold γ^t

$$\mathcal{L}_e = \sum_{i=1}^B H(C + 1, p_{\theta}(y[\Pi(\tilde{x}_i)])) \quad \text{Strong Augmentation}$$

$$\mathcal{L} = \mathcal{L}_s + \lambda(\mathcal{L}_u + \mathcal{L}_e)$$

Algorithm 1 MarginMatch

Require: Labeled data L ; unlabeled data U ; erroneous examples E ; maximum number of iterations T ; number of classes $C + 1$ (C original classes plus one virtual class of erroneous examples); θ model; π weak augmentations; Π strong augmentations.

- 1: Initialize the Average Pseudo-Margin (APM) threshold γ^1 at the first iteration to a small value (e.g., $\gamma^1 = -\infty$).
 - 2: **for** $t = 1$ to T **do**
 - 3: Estimate learning status α_c (using Eq. 2) and calculate the class-wise flexible thresholds \mathcal{T}_c^t (using Eq. 3) for each class c .
 - 4: **while** U not exhausted **do**
 - 5: Labeled batch $L_b = \{(x_1, y_1), \dots, (x_B, y_B)\}$, unlabeled batch $U_b = \{\hat{x}_1, \dots, \hat{x}_{\nu B}\}$, erroneous (or mislabeled) batch $E_b = \{(\tilde{x}_1, C + 1), \dots, (\tilde{x}_B, C + 1)\}$
 - 6: **for** $x \in U_b \cup E_b$ **do**
 - 7: Compute logits z_c for each class c after applying weak augmentations when $x \in U_b$ and strong augmentations when $x \in E_b$.
 - 8: Calculate pseudo-margin PM_c^t (using Eq. 5) and update Average PM_c^t (using Eq. 6) for each $c = 1$ to $C + 1$.
 - 9: **end for**
 - 10: Minimize $\mathcal{L} = \mathcal{L}_s + \lambda(\mathcal{L}_u + \mathcal{L}_e)$
 - 11: $\mathcal{L}_s = \frac{1}{B} \sum_{i=1}^B H(y_i, p_\theta(y|\pi(x_i)))$
 - 12: $\mathcal{L}_u = \sum_{i=1}^{\nu B} \mathbb{1}(\text{AM}_{\hat{p}_\theta(y|\pi(\hat{x}_i))}^t(\hat{x}_i) > \gamma^t) \times \mathbb{1}(\max(p_\theta(y|\pi(\hat{x}_i))) > \mathcal{T}_{\hat{p}_\theta(y|\pi(\hat{x}_i))}^t) \times H(\hat{p}_\theta(y|\pi(\hat{x}_i)), p_\theta(y|\Pi(\hat{x}_i)))$
 - 13: $\mathcal{L}_e = \sum_{i=1}^B H(C + 1, p_\theta(y|\Pi(\tilde{x}_i)))$
 - 14: **end while**
 - 15: Update γ^{t+1} as the 95th percentile erroneous sample APM_{C+1}^t .
 - 16: **end for**
-

Experiment

Dataset	CIFAR-10			CIFAR-100			SVHN			STL-10		
	4	25	400	4	25	100	4	25	100	4	25	100
Pseudo-Labeling	74.61 _{0.26}	46.49 _{2.20}	15.08 _{0.19}	87.45 _{0.85}	57.74 _{0.28}	36.55 _{0.24}	64.61 _{5.60}	25.21 _{2.03}	9.40 _{0.32}	74.68 _{0.99}	55.45 _{2.43}	32.64 _{0.71}
UDA	10.79 _{3.75}	5.32 _{0.06}	4.41 _{0.07}	48.95 _{1.59}	29.43 _{0.21}	23.87 _{0.23}	5.34 _{4.27}	4.26 _{0.39}	1.95 _{0.01}	37.82 _{8.44}	9.81 _{1.15}	6.81 _{0.17}
MixMatch	45.24 _{2.15}	12.76 _{1.14}	7.13 _{0.34}	62.15 _{2.17}	41.51 _{1.19}	28.16 _{0.24}	46.18 _{1.78}	3.98 _{0.17}	3.5 _{0.13}	34.15 _{1.54}	8.95 _{0.32}	10.41 _{0.73}
ReMixMatch	5.27 _{0.19}	4.85 _{0.13}	4.04 _{0.12}	47.15 _{0.76}	27.14 _{0.23}	23.78 _{0.12}	4.23 _{0.31}	3.18 _{0.04}	1.94 _{0.06}	31.51 _{0.75}	8.54 _{0.48}	6.19 _{0.24}
FixMatch	7.8 _{0.28}	4.91 _{0.05}	4.25 _{0.08}	48.21 _{0.82}	29.45 _{0.16}	22.89 _{0.12}	3.97 _{1.18}	3.13_{1.03}	1.97 _{0.03}	38.43 _{4.14}	10.45 _{1.04}	6.43 _{0.33}
FlexMatch	5.04 _{0.06}	5.04 _{0.09}	4.19 _{0.01}	39.99 _{1.62}	26.96 _{0.08}	22.44 _{0.15}	8.19 _{3.20}	7.78 _{2.55}	6.72 _{0.30}	29.15 _{1.32}	8.23 _{0.13}	5.77 _{0.12}
MarginMatch	4.91_{0.07}	4.73_{0.12}	3.98_{0.02}	36.97_{1.32}	23.71_{0.13}	21.39_{0.12}	3.75_{1.20}	3.14 _{1.17}	1.93_{0.01}	25.37_{3.58}	7.31_{0.35}	5.52_{0.15}

Table 1. Test error rates on CIFAR-10, CIFAR-100, SVHN, and STL-10 datasets. Best results are shown in **blue**.

Experiment

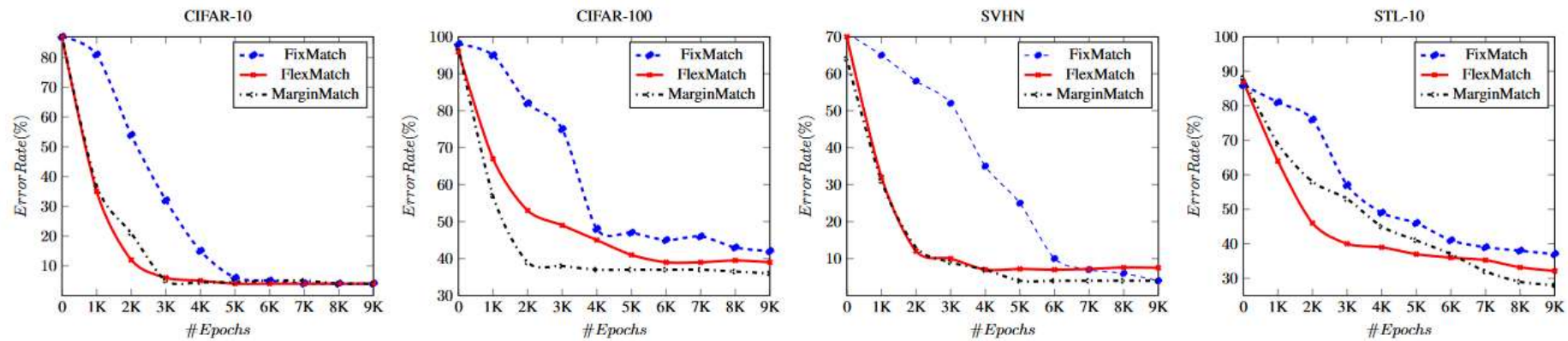


Figure 2. Convergence speed of MarginMatch against FixMatch and FlexMatch with 4 labels per class.

Dataset	ImageNet		WebVision	
	TOP-1	TOP-5	TOP-1	TOP-5
Supervised	48.39	25.49	49.58	26.78
FixMatch	43.66	21.80	44.76	22.65
FlexMatch	42.02	19.49	43.87	22.07
MarginMatch	41.05	18.28	43.08	21.13

Table 2. Test error rates on the ImageNet and WebVision datasets. Best results are shown in blue.

Experiment

Ablation Study

δ	0.95	0.99	0.995	0.997	0.999	1
ERR RATE	38.13	38.05	37.92	37.91	39.12	39.72

EMA coefficient

Table 3. Error rates obtained on CIFAR-100 with four examples per class and various smoothing values δ . Best result is in **blue**.

Dataset	CIFAR-10			CIFAR-100			SVHN			STL-10		
	4	25	400	4	25	100	4	25	100	4	25	100
Avg Confidence	23.87 _{2.73}	14.21 _{1.37}	7.54 _{0.78}	41.23 _{2.15}	31.49 _{1.48}	24.11 _{2.36}	8.99 _{4.27}	6.54 _{0.39}	4.73 _{0.01}	31.67 _{8.44}	14.87 _{1.15}	7.59 _{0.17}
Avg Entropy	8.58 _{0.41}	6.18 _{0.15}	5.85 _{0.12}	45.10 _{0.91}	26.02 _{1.11}	22.13 _{0.25}	15.69 _{1.25}	12.74 _{0.78}	9.33 _{0.05}	29.54 _{3.51}	10.63 _{1.35}	10.84 _{0.47}
Avg Margin	7.25 _{0.29}	5.38 _{0.76}	4.73 _{0.09}	39.72 _{1.52}	25.21 _{0.52}	23.18 _{0.17}	18.45 _{1.36}	11.29 _{0.93}	8.40 _{0.04}	28.45 _{4.28}	9.34 _{1.34}	7.59 _{0.21}
EMA Confidence	4.91 _{0.45}	4.74 _{0.09}	3.99 _{0.06}	38.67 _{0.74}	25.61 _{0.12}	21.48 _{0.17}	3.84 _{0.23}	3.25 _{0.03}	1.93 _{0.09}	25.9 _{0.81}	7.6 _{0.42}	5.74 _{0.57}
EMA Entropy	6.4 _{0.43}	8.34 _{0.12}	4.21 _{0.09}	41.63 _{0.76}	36.84 _{0.13}	22.52 _{0.07}	3.81 _{1.26}	3.17 _{0.87}	2.14 _{0.04}	27.21 _{4.05}	8.28 _{1.01}	6.79 _{0.27}
EMA Margin	4.91 _{0.07}	4.73 _{0.12}	3.98 _{0.02}	36.97 _{1.32}	23.71 _{0.13}	21.39 _{0.12}	3.75 _{1.20}	3.14 _{1.17}	1.93 _{0.01}	25.37 _{3.58}	7.31 _{0.35}	5.52 _{0.15}

Table 4. Test error rates comparing pseudo-margin with confidence and entropy. Best results are shown in **blue**.

$$\mathcal{L}_u = \sum_{i=1}^{\nu B} \mathbf{1} \left\{ \underbrace{\text{APM(AUM)}}_{\text{APM(AUM)}} \left(\text{AM}_{\hat{p}_\theta(y|\pi(\hat{x}_i))}^t(\hat{x}_i) > \gamma^t \right) \mathbf{1} \left(\underbrace{\text{FlexMatch}}_{\text{FlexMatch}} \left(\max(p_\theta(y|\pi(\hat{x}_i))) > \mathcal{T}_{\hat{p}_\theta(y|\pi(\hat{x}_i))}^t H(\hat{p}_\theta(y|\pi(\hat{x}_i)), p_\theta(y|\Pi(\hat{x}_i))) \right) \right\}$$

Experiment

Analysis

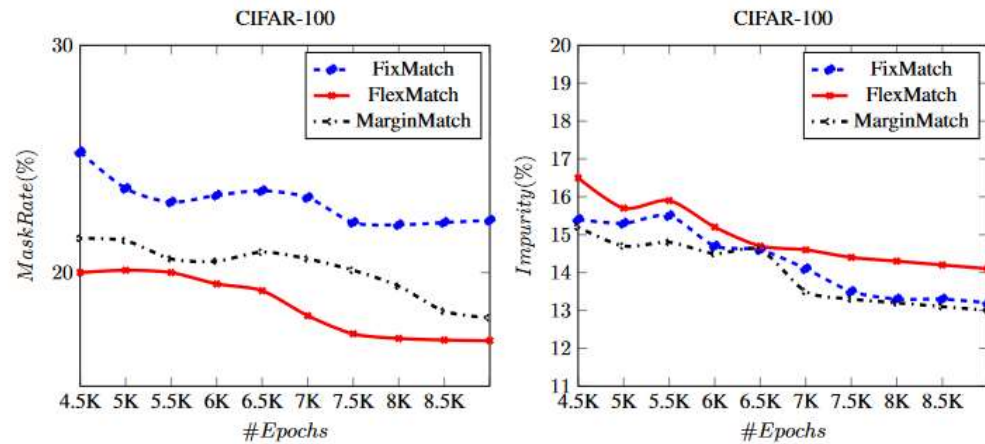


Figure 3. Mask rate and impurity on CIFAR-100 with 4 labeled examples per class.

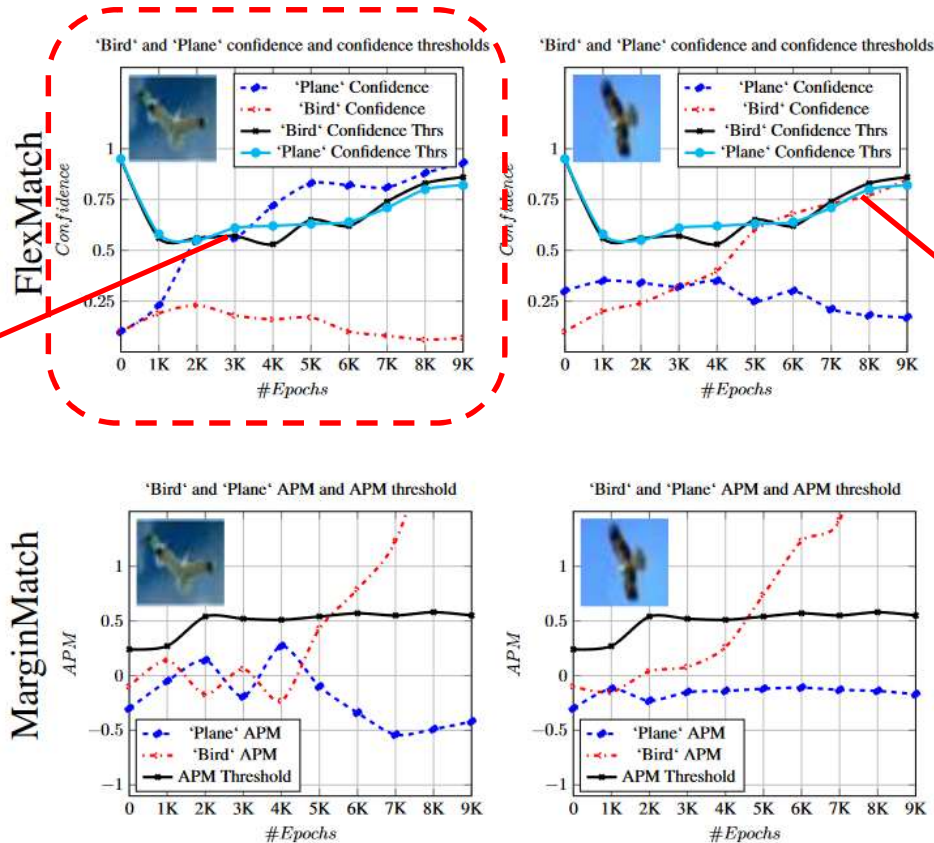
Mask rate is defined as the fraction of pseudo-labeled examples that do **not participate in** the training at epoch t due to confidence masking or pseudo-margin masking (or both).

Impurity in contrast is defined as the fraction of pseudolabeled examples that do **participate in** the training at epoch t **but with a wrong label**.

Experiment

Analysis

Oops, FlexMatch made a **mistake** here. Eventually this sample is wrongly classified.



FlexMatch is still **not** quite **confident** of its prediction at the end of the training.

Two bird images from CIFAR-10. Because of the **resemble** characteristics(background) between birds and plane, the model will get **confused**.