# FixMatch: Simplifying Semi-Supervised Learning with Consistency and Confidence<sup>1</sup>

February 6, 2020

 $<sup>^1{\</sup>it FixMatch:}$  Simplifying Semi-Supervised Learning with Consistency and Confidence by Sohn et al.

#### FixMatch:

▶ is a strong Self-Supervised Learning approach

- ▶ is a strong Self-Supervised Learning approach
- has a very simple loss function which has an interesting interpretation

- ▶ is a strong Self-Supervised Learning approach
- has a very simple loss function which has an interesting interpretation
- has quite complex augmentation strategy

- ▶ is a strong Self-Supervised Learning approach
- has a very simple loss function which has an interesting interpretation
- has quite complex augmentation strategy
- ▶ achieves 88.61% accuracy on CIFAR-10 with 4 labels per class!

▶ Imagine we have two datasets of images:

- ▶ Imagine we have two datasets of images:
  - ▶ Labeled images  $X^{I} = \{(x_{n}^{I}, p_{n})\}_{n=1}^{N}$

- ► Imagine we have two datasets of images:
  - ▶ Labeled images  $X' = \{(x_n', p_n)\}_{n=1}^N$
  - ▶ Unlabeled images  $X^u = \{x_n^u\}_{n=1}^{k \times N}$ , i.e.  $X^u$  is k times larger than  $X^l$

- ▶ Imagine we have two datasets of images:
  - ▶ Labeled images  $X^{I} = \{(x_{n}^{I}, p_{n})\}_{n=1}^{N}$
  - Unlabeled images  $X^u = \{x_n^u\}_{n=1}^{k \times N}$ , i.e.  $X^u$  is k times larger than  $X^l$
  - $\triangleright$   $p_n$  is a one-hot class label

- ▶ Imagine we have two datasets of images:
  - ▶ Labeled images  $X^{I} = \{(x_{n}^{I}, p_{n})\}_{n=1}^{N}$
  - Unlabeled images  $X^u = \{x_n^u\}_{n=1}^{k \times N}$ , i.e.  $X^u$  is k times larger than  $X^l$
  - $ightharpoonup p_n$  is a one-hot class label
- Let H(p, q) be a cross-entropy between p and q.

- ▶ Imagine we have two datasets of images:
  - ▶ Labeled images  $X^{I} = \{(x_{n}^{I}, p_{n})\}_{n=1}^{N}$
  - Unlabeled images  $X^u = \{x_n^u\}_{n=1}^{k \times N}$ , i.e.  $X^u$  is k times larger than  $X^l$
  - $ightharpoonup p_n$  is a one-hot class label
- Let H(p, q) be a cross-entropy between p and q.
- We also have augmentation functions:

- Imagine we have two datasets of images:
  - ▶ Labeled images  $X^{l} = \{(x_{n}^{l}, p_{n})\}_{n=1}^{N}$
  - Unlabeled images  $X^u = \{x_n^u\}_{n=1}^{k \times N}$ , i.e.  $X^u$  is k times larger than  $X^l$
  - $ightharpoonup p_n$  is a one-hot class label
- Let H(p,q) be a cross-entropy between p and q.
- ▶ We also have augmentation functions:
  - $\alpha(x)$  is a weak (i.e. simple) augmentation: random horizontal flipping and translations

- Imagine we have two datasets of images:
  - ▶ Labeled images  $X^{I} = \{(x_{n}^{I}, p_{n})\}_{n=1}^{N}$
  - Unlabeled images  $X^u = \{x_n^u\}_{n=1}^{k \times N}$ , i.e.  $X^u$  is k times larger than  $X^l$
  - p<sub>n</sub> is a one-hot class label
- Let H(p,q) be a cross-entropy between p and q.
- ▶ We also have augmentation functions:
  - $\alpha(x)$  is a weak (i.e. simple) augmentation: random horizontal flipping and translations
  - $\rightarrow$  A(x) is a strong (i.e. sophisticated) augmentation: color inversion, translation, contrast adjustment, etc

FixMatch loss is a combination of two losses:

$$\mathcal{L}_{\mathsf{FM}} = \mathcal{L}_{\mathsf{cls}} + \lambda_{\mathsf{pl}} \mathcal{L}_{\mathsf{pl}} \tag{1}$$

FixMatch loss is a combination of two losses:

$$\mathcal{L}_{\mathsf{FM}} = \mathcal{L}_{\mathsf{cls}} + \lambda_{\mathsf{pl}} \mathcal{L}_{\mathsf{pl}} \tag{1}$$

 $ightharpoonup \mathcal{L}_{cls}$  is a usual cross-entropy classification loss on  $X^I$  dataset.

FixMatch loss is a combination of two losses:

$$\mathcal{L}_{\mathsf{FM}} = \mathcal{L}_{\mathsf{cls}} + \lambda_{\mathsf{pl}} \mathcal{L}_{\mathsf{pl}} \tag{1}$$

- $ightharpoonup \mathcal{L}_{cls}$  is a usual cross-entropy classification loss on  $X^I$  dataset.
- $ightharpoonup \mathcal{L}_{pl}$  is a *pseudo-labeling loss*, i.e. a cross-entropy loss that uses synthetic targets produced by our model

Pseudo-labeling loss  $\mathcal{L}_{pl}$  equals:

$$\mathcal{L}_{\mathsf{pl}} = \frac{1}{\mathsf{k} \mathsf{N}} \sum_{n=1}^{\mathsf{k} \mathsf{N}} 1 \left[ \mathsf{max} \left( \bar{q}_{\mathsf{n}} \right) \ge \tau \right] H \left( \bar{q}_{\mathsf{n}}', p_{\mathsf{m}} \left( y | \mathcal{A} \left( x_{\mathsf{n}} \right) \right) \right) \tag{2}$$

#### Where:

- $lack q_n = p_m(y|\alpha(x_n)),$  i.e. class probabilities for weakly augmented  $x_n$
- $ightharpoonup ar{q}_n' = rg \max ar{q}_n$ , i.e. a pseudo label (one-hot)
- ightharpoonup au is a hyperparameter

Pseudo-labeling loss  $\mathcal{L}_{pl}$  equals:

$$\mathcal{L}_{\mathsf{pl}} = \frac{1}{\mathsf{k} \mathsf{N}} \sum_{n=1}^{\mathsf{k} \mathsf{N}} 1 \left[ \mathsf{max} \left( \bar{q}_{\mathsf{n}} \right) \ge \tau \right] H \left( \bar{q}_{\mathsf{n}}', p_{\mathsf{m}} \left( y | \mathcal{A} \left( x_{\mathsf{n}} \right) \right) \right) \tag{2}$$

Where:

- $\bar{q}_n = p_m(y|\alpha(x_n))$ , i.e. class probabilities for weakly augmented  $x_n$
- $ightharpoonup ar{q}_n' = rg \max ar{q}_n$ , i.e. a pseudo label (one-hot)
- ightharpoonup is a hyperparameter

Pseudo-labeling loss  $\mathcal{L}_{pl}$  equals:

$$\mathcal{L}_{\mathsf{pl}} = \frac{1}{\mathsf{k} \mathsf{N}} \sum_{n=1}^{\mathsf{k} \mathsf{N}} 1 \left[ \mathsf{max} \left( \bar{q}_{\mathsf{n}} \right) \ge \tau \right] H \left( \bar{q}_{\mathsf{n}}', p_{\mathsf{m}} \left( y | \mathcal{A} \left( x_{\mathsf{n}} \right) \right) \right) \tag{2}$$

#### Where:

- $\bar{q}_n = p_m(y|\alpha(x_n))$ , i.e. class probabilities for weakly augmented  $x_n$
- $\bar{q}'_n = \arg\max \bar{q}_n$ , i.e. a pseudo label (one-hot)
- ightharpoonup is a hyperparameter

#### Algorithm:

 $\triangleright$  Pick an unlabeled image  $x_n$ , produce two augmentated versions:

$$\bar{x}_n = \alpha(x)_n$$
 and  $\tilde{x}_n = \mathcal{A}(x_n)$ 

Pseudo-labeling loss  $\mathcal{L}_{pl}$  equals:

$$\mathcal{L}_{\mathsf{pl}} = \frac{1}{\mathsf{k} \mathsf{N}} \sum_{n=1}^{\mathsf{k} \mathsf{N}} 1 \left[ \mathsf{max} \left( \bar{q}_{\mathsf{n}} \right) \ge \tau \right] H \left( \bar{q}_{\mathsf{n}}', p_{\mathsf{m}} \left( y | \mathcal{A} \left( x_{\mathsf{n}} \right) \right) \right) \tag{2}$$

#### Where:

- $ightharpoonup ar{q}_n = p_m(y|\alpha(x_n))$ , i.e. class probabilities for weakly augmented  $x_n$
- $\bar{q}'_n = \arg\max \bar{q}_n$ , i.e. a pseudo label (one-hot)
- $\triangleright \tau$  is a hyperparameter

- Pick an unlabeled image  $x_n$ , produce two augmentated versions:  $\bar{x}_n = \alpha(x)_n$  and  $\tilde{x}_n = \mathcal{A}(x_n)$
- ► Compute class probabilities  $\bar{q}_n$  and  $\tilde{q}_n$  for  $\bar{x}_n$  and  $\tilde{x}_n$

Pseudo-labeling loss  $\mathcal{L}_{pl}$  equals:

$$\mathcal{L}_{\mathsf{pl}} = \frac{1}{\mathsf{k} \mathsf{N}} \sum_{n=1}^{\mathsf{k} \mathsf{N}} 1 \left[ \mathsf{max} \left( \bar{q}_{\mathsf{n}} \right) \ge \tau \right] H \left( \bar{q}'_{\mathsf{n}}, \rho_{\mathsf{m}} \left( y | \mathcal{A} \left( x_{\mathsf{n}} \right) \right) \right) \tag{2}$$

#### Where:

- $\bar{q}_n = p_m(y|\alpha(x_n))$ , i.e. class probabilities for weakly augmented  $x_n$
- $ightharpoonup ar{q}'_n = \arg\maxar{q}_n$ , i.e. a pseudo label (one-hot)
- ightharpoonup is a hyperparameter

- Pick an unlabeled image  $x_n$ , produce two augmentated versions:  $\bar{x}_n = \alpha(x)_n$  and  $\tilde{x}_n = \mathcal{A}(x_n)$
- ► Compute class probabilities  $\bar{q}_n$  and  $\tilde{q}_n$  for  $\bar{x}_n$  and  $\tilde{x}_n$
- Pick only examples with confident class probabilities for weakly-augmented images

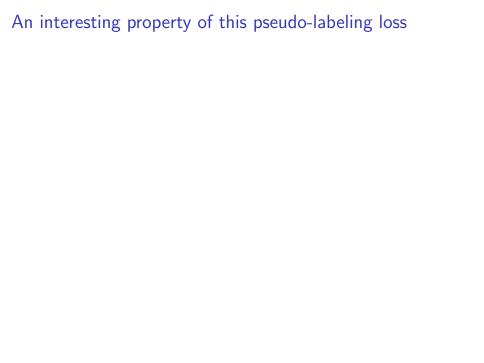
Pseudo-labeling loss  $\mathcal{L}_{pl}$  equals:

$$\mathcal{L}_{\mathsf{pl}} = \frac{1}{\mathsf{k} \mathsf{N}} \sum_{n=1}^{\mathsf{k} \mathsf{N}} 1 \left[ \mathsf{max} \left( \bar{q}_{\mathsf{n}} \right) \ge \tau \right] H \left( \bar{q}_{\mathsf{n}}', p_{\mathsf{m}} \left( y | \mathcal{A} \left( x_{\mathsf{n}} \right) \right) \right) \tag{2}$$

#### Where:

- $\bar{q}_n = p_m(y|\alpha(x_n))$ , i.e. class probabilities for weakly augmented  $x_n$
- $ightharpoonup ar{q}_n' = rg \max ar{q}_n$ , i.e. a pseudo label (one-hot)
- ightharpoonup au is a hyperparameter

- Pick an unlabeled image  $x_n$ , produce two augmentated versions:  $\bar{x}_n = \alpha(x)_n$  and  $\tilde{x}_n = \mathcal{A}(x_n)$
- ► Compute class probabilities  $\bar{q}_n$  and  $\tilde{q}_n$  for  $\bar{x}_n$  and  $\tilde{x}_n$
- Pick only examples with confident class probabilities for weakly-augmented images
- Cross-entropy term forces the model to give the same predictions for a weakly-augmented and a strongly-augmented images



# An interesting property of this pseudo-labeling loss

Previously used variants of pseudo-labelling loss required tuning of  $\lambda_{\rm pl}$  weight during training and gradually increase it.

# An interesting property of this pseudo-labeling loss

- Previously used variants of pseudo-labelling loss required tuning of  $\lambda_{\rm pl}$  weight during training and gradually increase it.
- In FixMatch model becomes gradually more confident in new images and authors omit  $\lambda_{pl}$  whatsoever!

# An interesting property of this pseudo-labeling loss

- Previously used variants of pseudo-labelling loss required tuning of  $\lambda_{\rm pl}$  weight during training and gradually increase it.
- In FixMatch model becomes gradually more confident in new images and authors omit  $\lambda_{\rm pl}$  whatsoever!
- So we get a curriculum learning out of the box!

► FixMatch provides very strong scores

- ► FixMatch provides very strong scores
- ► Has a very simple loss with an interesting side-effect of curriculum learning

- ► FixMatch provides very strong scores
- ► Has a very simple loss with an interesting side-effect of curriculum learning
- ► A disadvantage: strong augmentations are based on CutOut, CTAugment, etc and seem *very* sophisticated

- ► FixMatch provides very strong scores
- ► Has a very simple loss with an interesting side-effect of curriculum learning
- ► A disadvantage: strong augmentations are based on CutOut, CTAugment, etc and seem *very* sophisticated
- Strongly beats SotA in many setups:

Method	CIFAR-10			CIFAR-100			SVHN		
	40 labels	250 labels	4000 labels	400 labels	2500 labels	10000 labels	40 labels	250 labels	1000 labels
Π-Model	-	54.26±3.97	14.01±0.38	-	57.25±0.48	37.88±0.11		18.96±1.92	7.54±0.36
Pseudo-Labeling	-	$49.78 \pm 0.43$	$16.09 \pm 0.28$	-	$57.38\pm0.46$	$36.21\pm0.19$	-	$20.21\pm1.09$	$9.94\pm0.61$
Mean Teacher	-	$32.32\pm2.30$	$9.19\pm0.19$	-	53.91±0.57	$35.83 \pm 0.24$	-	$3.57\pm0.11$	$3.42\pm0.07$
MixMatch	47.54±11.50	$11.05\pm0.86$	$6.42\pm0.10$	$67.61\pm1.32$	$39.94\pm0.37$	$28.31 \pm 0.33$	$42.55\pm14.53$	$3.98\pm0.23$	$3.50\pm0.28$
UDA	29.05±5.93	$8.82 \pm 1.08$	$4.88\pm0.18$	$59.28 \pm 0.88$	$33.13 \pm 0.22$	$24.50\pm0.25$	$52.63\pm20.51$	$5.69 \pm 2.76$	2.46±0.24
ReMixMatch	19.10±9.64	5.44±0.05	$4.72\pm0.13$	$44.28 \pm 2.06$	27.43±0.31	$23.03 \pm 0.56$	$3.34 \pm 0.20$	$2.92 \pm 0.48$	$2.65 \pm 0.08$
FixMatch (RA)	13.81±3.37	5.07±0.65	4.26±0.05	48.85±1.75	28.29±0.11	22.60±0.12	3.96±2.17	2.48±0.38	2.28±0.11
FixMatch (CTA)	$11.39 \pm 3.35$	$5.07 \pm 0.33$	$4.31\pm0.15$	49.95±3.01	$28.64 \pm 0.24$	$23.18\pm0.11$	$7.65\pm7.65$	$2.64\pm0.64$	2.36±0.19