# REMIXMATCH: Semi-supervised learning with distribution alignment and augmentation anchoring [ICLR 2020]

Distribution Alignment: the marginal distribution of predictions on $D_u$ to be close to ———————————————— ground truth labels

Augmentation Anchoring: Feed mulitiple "CTAugment" Strongly augmented versions of an input encourage each output to be closed to the weakly augmentation.

## MixMatch:

$$\chi = \{(x_b, p_b) : b \in (1, \cdots, B)\}$$

$$\mathcal{U} = \{u_b : b \in (1, \cdots, B)\} \overset{u_b}{\underset{\vdots}{\rightrightarrows}} \hat{u}_{b,K}, K \in \{1, \cdots, K\} \xrightarrow{\text{average}} \bar{q}_b = \frac{1}{K}\sum_K P(y|\hat{u}_{b,k}, \theta)$$

K weakly augmentation

$\downarrow$ Sharpen

guess labels

$$\mathcal{U}_g = \{(u_b, q_b)\}$$

Combine $\chi$ and $\mathcal{U}_g \rightarrow \mathcal{y}$

MixUp: $(x', p') = \lambda(x_1, p_1) + (1-\lambda)(x_2, p_2), \forall (x_1, p_1), (x_2, p_2) \in \mathcal{y}$

Given these mixed-up samples, it performs standrad fully-supervised training with minor modifications.

ReMixMatch:

① Distribution Alignment

Input-Output mutual information: (maximize)

$$I(y; x) = \iint p(y, x) \log \frac{P(y, x)}{P(y) \, P(x)} \, dy \, dx$$

$$= \mathcal{H}\left(\mathbb{E}_x\left[P_{model}(y|x; \theta)\right]\right) - \mathbb{E}_x\left[\mathcal{H}\left(P_{model}(y|x; \theta)\right)\right]$$

<span style="color:blue">dataset's marginal class</span> <span style="color:red">"fair"</span>
<span style="color:blue">distribution $p(y)$ uniform X</span>  ↓

<span style="color:red">entropy minimization
(high confidence in a class label)</span>

maintain a running average of the model's prediction on $U$
given the model's prediction $\underline{q = P_{model}(y|u; \theta)}$   $\tilde{P}(y)$

Scale $q$ by a ratio and renormalize the result

$$\tilde{q} = Normalize\left(q \times \frac{p(y)}{\tilde{p}(y)}\right)$$
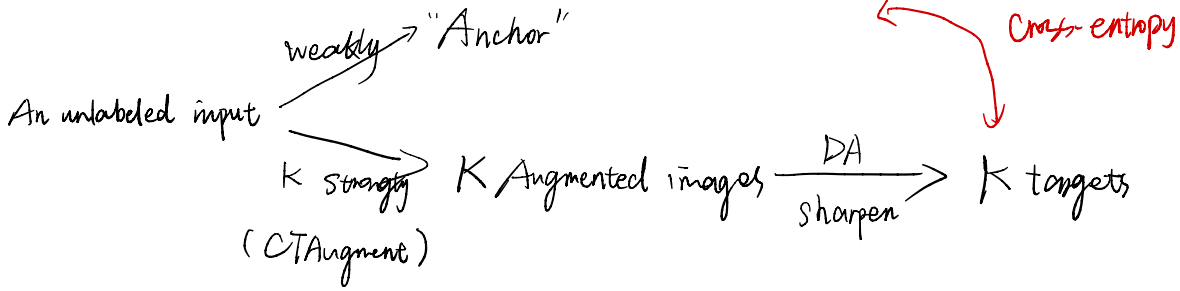
<span style="color:blue">eg. $q = (\frac{1}{2}, \frac{1}{4}, \frac{1}{4})$, $\tilde{p}(y) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$, $p(y) = (1, 0, 0)$</span>

<span style="color:blue">$\tilde{q} = Normalize\left((\frac{1}{2}, \frac{1}{4}, \frac{1}{4}) \times (\frac{1}{3}, 0, 0)\right) = (1, 0, 0)$</span>

Note that $p(y)$ is estimated by the labeled examples seen during training.
Or it is known a priori.

② Augmentation Anchoring

weakly → "Anchor"

An unlabeled input

← Cross-entropy

K strongly → K Augmented images $\xrightarrow[\text{sharpen}]{DA}$ → K targets

(CTAugment)

Control Theory Augment (CTAugment)

A method for learning a data augmentation policy which results in high validation set accuracy.

---

**Algorithm 1** ReMixMatch algorithm for producing a collection of processed labeled examples and processed unlabeled examples with label guesses (cf. Berthelot et al. (2019) Algorithm 1.)

---

1: **Input:** Batch of labeled examples and their one-hot labels $\mathcal{X} = \{(x_b, p_b) : b \in (1, \ldots, B)\}$, batch of unlabeled examples $\mathcal{U} = \{u_b : b \in (1, \ldots, B)\}$, sharpening temperature $T$, number of augmentations $K$, Beta distribution parameter $\alpha$ for MixUp.
2: **for** $b = 1$ **to** $B$ **do**
3:     $\hat{x}_b = \text{StrongAugment}(x_b)$     // Apply strong data augmentation to $x_b$
4:     $\hat{u}_{b,k} = \text{StrongAugment}(u_b); k \in \{1, \ldots, K\}$     // Apply strong data augmentation $K$ times to $u_b$
5:     $\tilde{u}_b = \text{WeakAugment}(u_b)$     // Apply weak data augmentation to $u_b$
6:     $q_b = p_{\text{model}}(y \mid \tilde{u}_b; \theta)$     // Compute prediction for weak augmentation of $u_b$
7:     $q_b = \text{Normalize}(q_b \times p(y)/\tilde{p}(y))$     // Apply distribution alignment
8:     $q_b = \text{Normalize}(q_b^{1/T})$     // Apply temperature sharpening to label guess
9: **end for**
10: $\hat{\mathcal{X}} = ((\hat{x}_b, p_b); b \in (1, \ldots, B))$     // Augmented labeled examples and their labels
11: $\hat{\mathcal{U}}_1 = ((\hat{u}_{b,1}, q_b); b \in (1, \ldots, B))$     // First strongly augmented unlabeled example and guessed label
12: $\hat{\mathcal{U}} = ((\hat{u}_{b,k}, q_b); b \in (1, \ldots, B), k \in (1, \ldots, K))$     // All strongly augmented unlabeled examples
13: $\hat{\mathcal{U}} = \hat{\mathcal{U}} \cup ((\tilde{u}_b, q_b); b \in (1, \ldots, B))$     // Add weakly augmented unlabeled examples
14: $\mathcal{W} = \text{Shuffle}(\text{Concat}(\hat{\mathcal{X}}, \hat{\mathcal{U}}))$     // Combine and shuffle labeled and unlabeled data
15: $\mathcal{X}' = (\text{MixUp}(\hat{\mathcal{X}}_i, \mathcal{W}_i); i \in (1, \ldots, |\hat{\mathcal{X}}|))$     // Apply MixUp to labeled data and entries from $\mathcal{W}$
16: $\mathcal{U}' = (\text{MixUp}(\hat{\mathcal{U}}_i, \mathcal{W}_{i+|\hat{\mathcal{X}}|}); i \in (1, \ldots, |\hat{\mathcal{U}}|))$     // Apply MixUp to unlabeled data and the rest of $\mathcal{W}$
17: **return** $\mathcal{X}', \mathcal{U}', \hat{\mathcal{U}}_1$

↳ two additional loss

**Pre-mixup unlabeled loss**    We feed the guessed labels and predictions for example in $\hat{\mathcal{U}}_1$ as-is into a separate cross-entropy loss term.

**Rotation loss**    Recent result have shown that applying ideas from self-supervised learning to SSL can produce strong performance (Gidaris et al., 2018; Zhai et al., 2019). We integrate this idea by rotating each image $u \in \hat{\mathcal{U}}_1$ as $\text{Rotate}(u, r)$ where we sample the rotation angle $r$ uniformly from $r \sim \{0, 90, 180, 270\}$ and then ask the model to predict the rotation amount as a four-class classification problem.

In total, the ReMixMatch loss is

$$\sum_{x,p \in \mathcal{X}'} \text{H}(p, p_{\text{model}}(y|x; \theta)) + \lambda_{\mathcal{U}} \sum_{u,q \in \mathcal{U}'} \text{H}(q, p_{\text{model}}(y|u; \theta)) \tag{3}$$

$$+ \lambda_{\hat{\mathcal{U}}_1} \sum_{u,q \in \hat{\mathcal{U}}_1} \text{H}(q, p_{\text{model}}(y|u; \theta)) + \lambda_r \sum_{u \in \hat{\mathcal{U}}_1} \text{H}(r, p_{\text{model}}(r| \text{Rotate}(u, r); \theta)) \tag{4}$$