



# MixMatch: A Holistic Approach to Semi-Supervised Learning

---

**David Berthelot**  
Google Research  
dberth@google.com

**Nicholas Carlini**  
Google Research  
ncarlini@google.com

**Ian Goodfellow**  
Work done at Google  
ian-academic@mailfence.com

**Avital Oliver**  
Google Research  
avitalo@google.com

**Nicolas Papernot**  
Google Research  
papernot@google.com

**Colin Raffel**  
Google Research  
craffel@google.com

**NIPS-2019**

# Related Work for SSL

---

## ■ Consistency Regularization

A classifier should output the same class distribution for an unlabeled example even after it has been augmented.

$$\|p_{\text{model}}(y \mid \text{Augment}(x); \theta) - p_{\text{model}}(y \mid \text{Augment}(x); \theta)\|_2^2.$$

## ■ Entropy Minimization

The classifier's decision boundary should not pass through high-density regions of the marginal data distribution.

- Require the classifier output low-entropy predictions on unlabeled data.

## ■ Traditional Regularization

Make it harder to memorize the training data and therefore hopefully make it generalize better to unseen data

- Penalizes the L2 norm of the model parameters
- MixUp

# MixMatch

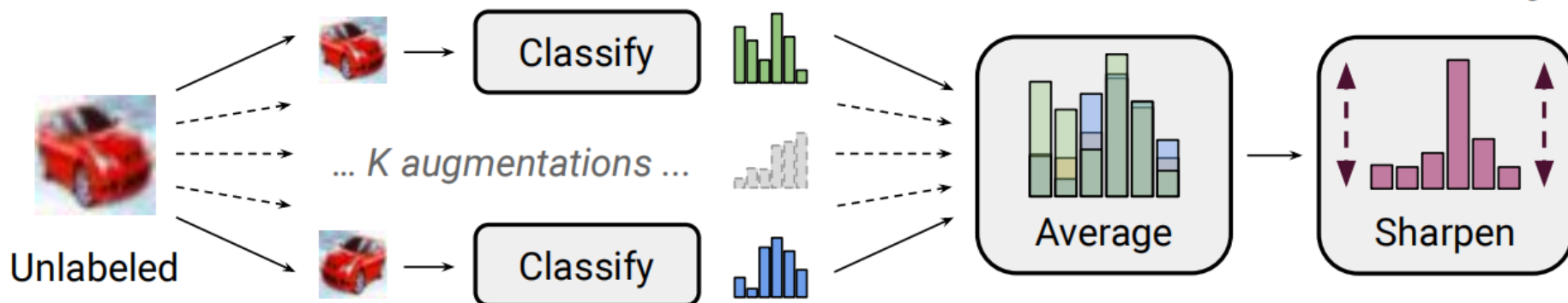
## ■ Data Augmentation

$$\hat{x}_b = \text{Augment}(x_b)$$

$$\hat{u}_{b,k} = \overline{\text{Augment}}(\overline{u}_b), k \in (1, \dots, K) \quad K \text{ augmentations}$$

## ■ Label Guessing

$$\bar{q}_b = \frac{1}{K} \sum_{k=1}^K P_{\text{model}}(y \mid \hat{u}_{b,k}; \theta) \quad \text{Sharpen}(p, T)_i := p_i^{\frac{1}{T}} / \sum_{j=1}^L p_j^{\frac{1}{T}}$$



# MixMatch

## ■ MixUp

Mix both labeled examples and unlabeled examples with label guesses  
Corresponding labels probabilities  $(x_1, p_1), (x_2, p_2)$  we compute  $(x', p')$

$$\lambda \sim \text{Beta}(\alpha, \alpha)$$

$$\lambda' = \max(\lambda, 1 - \lambda)$$

$$x' = \lambda' x_1 + (1 - \lambda') x_2$$

$$p' = \lambda' p_1 + (1 - \lambda') p_2$$

Linear interpolations of feature vectors should lead to linear interpolations of the associated targets.

Encourages the model  $f$  to behave linearly in-between training examples.

Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz.

mixup: Beyond empirical risk minimization. ICLR-2018

# MixMatch

## ■ MixUp

Mix both labeled examples and unlabeled examples with label guesses  
Corresponding labels probabilities  $(x_1, p_1), (x_2, p_2)$  we compute  $(x', p')$

$$\lambda \sim \text{Beta}(\alpha, \alpha)$$

$$\lambda' = \max(\lambda, 1 - \lambda)$$

$$x' = \lambda' x_1 + (1 - \lambda') x_2$$

$$p' = \lambda' p_1 + (1 - \lambda') p_2$$

$\hat{\mathcal{X}} = ((\hat{x}_b, p_b); b \in (1, \dots, B))$  // Augmented labeled examples and their labels

$\hat{\mathcal{U}} = ((\hat{u}_{b,k}, q_b); b \in (1, \dots, B), k \in (1, \dots, K))$  // Augmented unlabeled examples, guessed labels

$\mathcal{W} = \text{Shuffle}(\text{Concat}(\hat{\mathcal{X}}, \hat{\mathcal{U}}))$  // Combine and shuffle labeled and unlabeled data

$\mathcal{X}' = (\text{MixUp}(\hat{\mathcal{X}}_i, \mathcal{W}_i); i \in (1, \dots, |\hat{\mathcal{X}}|))$  // Apply MixUp to labeled data and entries from  $\mathcal{W}$

$\mathcal{U}' = (\text{MixUp}(\hat{\mathcal{U}}_i, \mathcal{W}_{i+|\hat{\mathcal{X}}|}); i \in (1, \dots, |\hat{\mathcal{U}}|))$  // Apply MixUp to unlabeled data and the rest of  $\mathcal{W}$

# MixMatch

## ■ Loss Function

$$\mathcal{X}', \mathcal{U}' = \text{MixMatch}(\mathcal{X}, \mathcal{U}, T, K, \alpha)$$

$$\mathcal{L}_{\mathcal{X}} = \frac{1}{|\mathcal{X}'|} \sum_{x, p \in \mathcal{X}'} \text{H}(p, \text{P}_{\text{model}}(y | x; \theta))$$

$$\mathcal{L}_{\mathcal{U}} = \frac{1}{L|\mathcal{U}'|} \sum_{u, q \in \mathcal{U}'} \|q - \text{P}_{\text{model}}(y | u; \theta)\|_2^2$$

$$\mathcal{L} = \mathcal{L}_{\mathcal{X}} + \lambda_{\mathcal{U}} \mathcal{L}_{\mathcal{U}}$$

## ■ Hyperparameters

$$T = 0.5 \quad K = 2 \quad \alpha = 0.75 \quad \lambda_{\mathcal{U}} = 100$$

# Algorithm of MixMatch

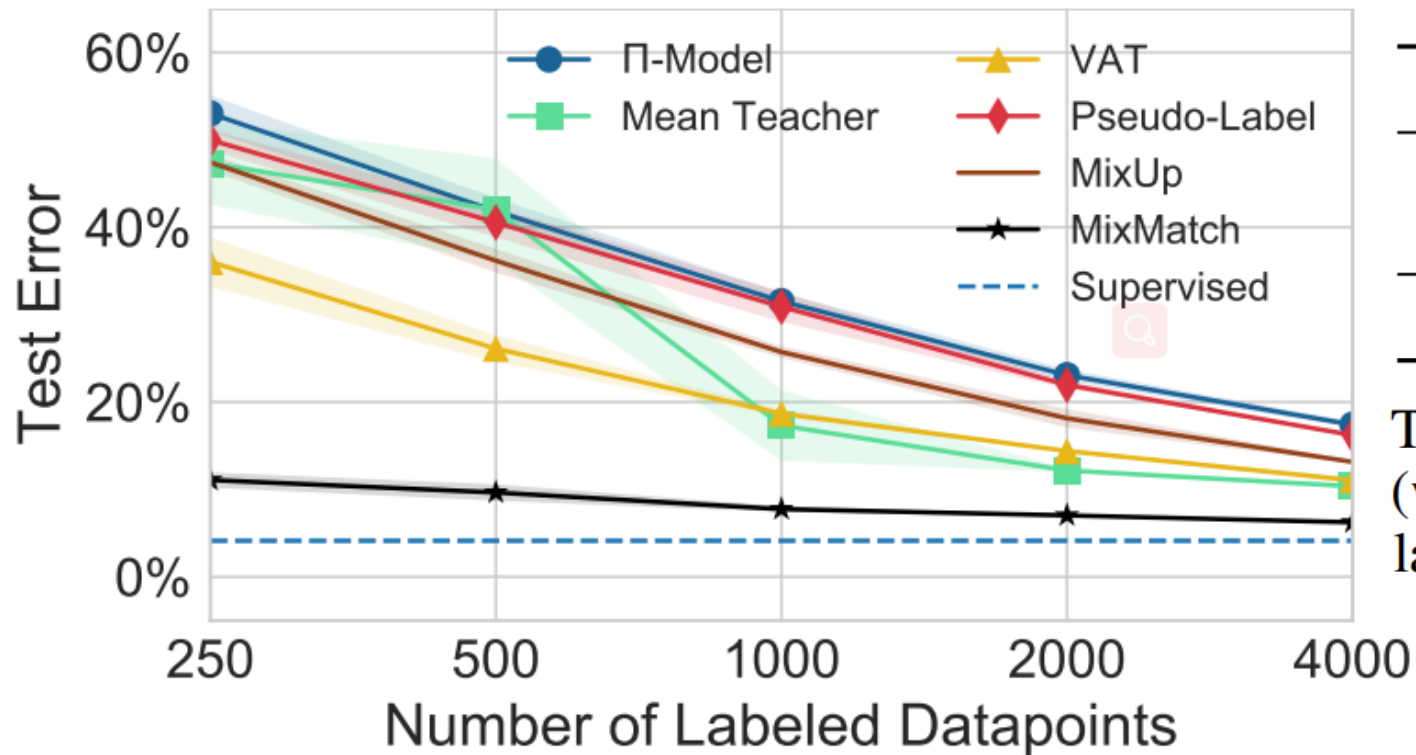
**Algorithm 1** MixMatch takes a batch of labeled data  $\mathcal{X}$  and a batch of unlabeled data  $\mathcal{U}$  and produces a collection  $\mathcal{X}'$  (resp.  $\mathcal{U}'$ ) of processed labeled examples (resp. unlabeled with guessed labels).

- 1: **Input:** Batch of labeled examples and their one-hot labels  $\mathcal{X} = ((x_b, p_b); b \in (1, \dots, B))$ , batch of unlabeled examples  $\mathcal{U} = (u_b; b \in (1, \dots, B))$ , sharpening temperature  $T$ , number of augmentations  $K$ , Beta distribution parameter  $\alpha$  for MixUp.
- 2: **for**  $b = 1$  **to**  $B$  **do**
- 3:    $\hat{x}_b = \text{Augment}(x_b)$    *// Apply data augmentation to  $x_b$*
- 4:   **for**  $k = 1$  **to**  $K$  **do**
- 5:      $\hat{u}_{b,k} = \text{Augment}(u_b)$    *// Apply  $k^{\text{th}}$  round of data augmentation to  $u_b$*
- 6:   **end for**
- 7:    $\bar{q}_b = \frac{1}{K} \sum_k \text{P}_{\text{model}}(y \mid \hat{u}_{b,k}; \theta)$    *// Compute average predictions across all augmentations of  $u_b$*
- 8:    $q_b = \text{Sharpen}(\bar{q}_b, T)$    *// Apply temperature sharpening to the average prediction (see eq. (7))*
- 9: **end for**
- 10:  $\hat{\mathcal{X}} = ((\hat{x}_b, p_b); b \in (1, \dots, B))$    *// Augmented labeled examples and their labels*
- 11:  $\hat{\mathcal{U}} = ((\hat{u}_{b,k}, q_b); b \in (1, \dots, B), k \in (1, \dots, K))$    *// Augmented unlabeled examples, guessed labels*
- 12:  $\mathcal{W} = \text{Shuffle}(\text{Concat}(\hat{\mathcal{X}}, \hat{\mathcal{U}}))$    *// Combine and shuffle labeled and unlabeled data*
- 13:  $\mathcal{X}' = (\text{MixUp}(\hat{\mathcal{X}}_i, \mathcal{W}_i); i \in (1, \dots, |\hat{\mathcal{X}}|))$    *// Apply MixUp to labeled data and entries from  $\mathcal{W}$*
- 14:  $\mathcal{U}' = (\text{MixUp}(\hat{\mathcal{U}}_i, \mathcal{W}_{i+|\hat{\mathcal{X}}|}); i \in (1, \dots, |\hat{\mathcal{U}}|))$    *// Apply MixUp to unlabeled data and the rest of  $\mathcal{W}$*
- 15: **return**  $\mathcal{X}', \mathcal{U}'$

# Experiment

■ Wide ResNet-28

$\lambda_U = 75$  for CIFAR-10 and  $\lambda_U = 150$  for CIFAR-100



Method	CIFAR-10	CIFAR-100
Mean Teacher [44]	6.28	-
SWA [2]	5.00	28.80
MixMatch	$4.95 \pm 0.08$	$25.88 \pm 0.30$

Table 1: CIFAR-10 and CIFAR-100 error rate (with 4,000 and 10,000 labels respectively) with larger models (26 million parameters).

# Experiment

## B.1 CIFAR-10

Training the same model with supervised learning on the entire 50000-example training set achieved an error rate of 4.13%.

Methods/Labels	250	500	1000	2000	4000
PiModel	$53.02 \pm 2.05$	$41.82 \pm 1.52$	$31.53 \pm 0.98$	$23.07 \pm 0.66$	$17.41 \pm 0.37$
PseudoLabel	$49.98 \pm 1.17$	$40.55 \pm 1.70$	$30.91 \pm 1.73$	$21.96 \pm 0.42$	$16.21 \pm 0.11$
Mixup	$47.43 \pm 0.92$	$36.17 \pm 1.36$	$25.72 \pm 0.66$	$18.14 \pm 1.06$	$13.15 \pm 0.20$
VAT	$36.03 \pm 2.82$	$26.11 \pm 1.52$	$18.68 \pm 0.40$	$14.40 \pm 0.15$	$11.05 \pm 0.31$
MeanTeacher	$47.32 \pm 4.71$	$42.01 \pm 5.86$	$17.32 \pm 4.00$	$12.17 \pm 0.22$	$10.36 \pm 0.25$
MixMatch	$11.08 \pm 0.87$	$9.65 \pm 0.94$	$7.75 \pm 0.32$	$7.03 \pm 0.15$	$6.24 \pm 0.06$

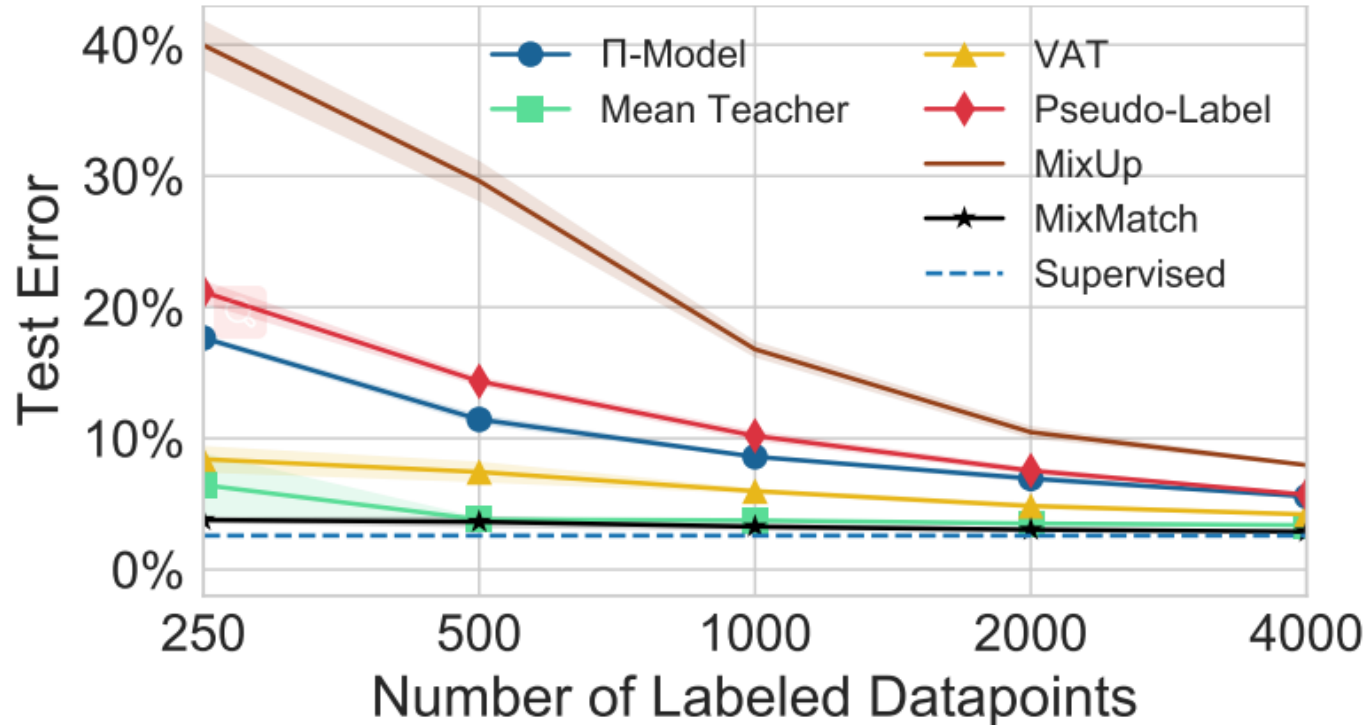
Table 5: Error rate (%) for CIFAR10.

# Experiment

■ Wide ResNet-28

$\lambda_U = 250$  for SVHN

SVHN has two training sets: train and extra



Labels	250	500	1000	2000	4000	All
SVHN	$3.78 \pm 0.26$	$3.64 \pm 0.46$	$3.27 \pm 0.31$	$3.04 \pm 0.13$	$2.89 \pm 0.06$	2.59
SVHN+Extra	$2.22 \pm 0.08$	$2.17 \pm 0.07$	$2.18 \pm 0.06$	$2.12 \pm 0.03$	$2.07 \pm 0.05$	1.71

Table 3: Comparison of error rates for SVHN and SVHN+Extra for MixMatch. The last column (“All”) contains the fully-supervised performance with all labels in the corresponding training set.

# Experiment

- Wide ResNet-28       $\lambda_U = 50$  for STL-10

Method	1000 labels	5000 labels
CutOut [12]	-	12.74
IIC [20]	-	11.20
SWWAE [48]	25.70	-
CC-GAN <sup>2</sup> [11]	22.20	-
MixMatch	10.18 $\pm$ 1.46	5.59

Table 2: STL-10 error rate using 1000-label splits or the entire 5000-label training set.

# Ablation Study

Ablation	250 labels	4000 labels
MixMatch	11.80	6.00
MixMatch without distribution averaging ( $K = 1$ )	17.09	8.06
MixMatch with $K = 3$	11.55	6.23
MixMatch with $K = 4$	12.45	5.88
MixMatch without temperature sharpening ( $T = 1$ )	27.83	10.59
MixMatch with parameter EMA	11.86	6.47
MixMatch without MixUp	39.11	10.97
MixMatch with MixUp on labeled only	32.16	9.22
MixMatch with MixUp on unlabeled only	12.35	6.83
MixMatch with MixUp on separate labeled and unlabeled	12.26	6.50
Interpolation Consistency Training [45]	38.60	6.81

Table 4: Ablation study results. All values are error rates on CIFAR-10 with 250 or 4000 labels.