**MixMatch: A Holistic Approach to Semi-Supervised Learning**

David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, Colin Raffel
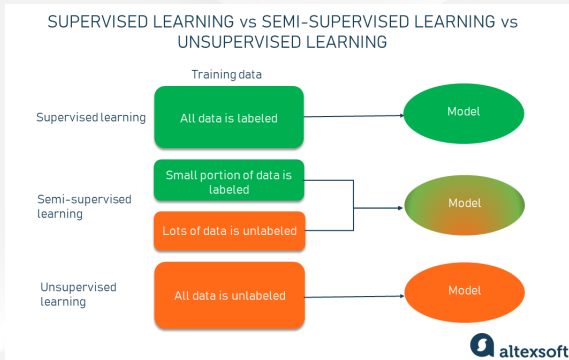
Coby Penso

Semi-supervised learning is an approach that combines a small amount of labeled data with a large amount of unlabeled data during training.



SUPERVISED LEARNING vs SEMI-SUPERVISED LEARNING vs UNSUPERVISED LEARNING

Training data

Supervised learning — All data is labeled → Model

Semi-supervised learning — Small portion of data is labeled / Lots of data is unlabeled → Model

Unsupervised learning — All data is unlabeled → Model

altexsoft

- data augmentations - applies input transformations assumed to leave class semantics unaffected
- **Consistency regularization** applies data augmentation to semi-supervised learning by leveraging the idea that a classifier should output the same class distribution for an unlabeled example even after it has been augmented. More formally, consistency regularization enforces that an unlabeled example x should be classified the same as Augment(x), an augmentation of itself.
- For unlabeled points x, Consistency Regularization loss term is

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

- Augment(x) - stochastic transformation
- **Drawback:** use domain specific data augmentations strategies.

## Entropy Minimization

- SSL assumes that the classifier's decision boundary should not pass through high-density regions of the marginal data distribution.
- One way to enforce - require that the classifier output low-entropy predictions on unlabeled data.
- Loss term is

$$p_{model}(y|x_{unlabeled}; \theta)$$

- "Pseudo-Label" does it implicitly.

Imposing a constraint on a model to make it harder to memorize the training data and therefore hopefully make it generalize better to unseen data.

- Weight decay
- L2 Regularization
- MixUp -
  to encourage convex behaviour "between" examples.

## MixMatch - Overview

- Batch *mathcalX* - labeled data.
- Batch $\mathcal{U}$ - equally-sized unlabeled data.
- Generate $\mathcal{X}'$ and $\mathcal{U}'$

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

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

$$\mathcal{L}_{\mathcal{U}} = \frac{1}{L|\mathcal{U}'|} \sum_{u,q \in \mathcal{U}'} H(q, p_{model}(y|u; \theta))$$

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

- H(p,q) - cross-entropy.
- T,K,*alpha*, and $\lambda_{\mathcal{U}}$ are hyperparameters.

- Both on labeled and unlabeled data.
- For each $x_b$ in $\mathcal{X}$: $\hat{x}_b = Augment(x_b)$.
- For each $u_b$ in $\mathcal{U}$: $\hat{u}_{b,k} = Augment(u_b)$, $k \in (1,..,K)$



Figure 1: Diagram of the label guessing process used in MixMatch. Stochastic data augmentation is applied to an unlabeled image $K$ times, and each augmented image is fed through the classifier. Then, the average of these $K$ predictions is "sharpened" by adjusting the distribution's temperature. See algorithm 1 for a full description.

- For each unlabeled example in $\mathcal{U}$, MixMatch produces a "guess" label.
- Later on, used in the unsupervised loss term.
- 
$$\hat{q}_b = \frac{1}{K} \sum_{k=1}^{K} p_{model}(y | \hat{u}_{b,k}; \theta)$$

- Apply a sharpening function to reduce the entropy of the label

$$Sharpen(p, T)_i := p_i^{\frac{1}{T}} / \sum_{j=1}^{L} p_j^{\frac{1}{T}}$$

Lower T -> lower-entropy predictions. distr

## MixMatch - Mixup

- Modified version of MixUp

$$(x_1, p_1), (x_2, p_2) \rightarrow (x', p')$$

$$\lambda \sim Beta(\alpha, \alpha)$$

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

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

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

Vanilla MixUp uses $\lambda' = \lambda$

- $\lambda' = max(\lambda, 1 - \lambda)$ ensures that x' is closer to $x_1$.
- MixMatch transforms $\mathcal{X}$ into $\mathcal{X}'$, a collection of labeled examples which have had data augmentation and MixUp (potentially mixed with an unlabeled example) applied.
- Similarly, $\mathcal{U}$ is transformed into $\mathcal{U}'$, a collection of multiple augmentations of each unlabeled example with corresponding label guesses.

**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 $\mathcal{X} = ((x_b, p_b); b \in (1, \ldots, B))$, 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:     $\hat{x}_b = \text{Augment}(x_b)$    // *Apply data augmentation to $x_b$*

4:     **for** $k = 1$ **to** $K$ **do**

5:        $\hat{u}_{b,k} = \text{Augment}(u_b)$    // *Apply $k^{th}$ 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: $\hat{\mathcal{X}} = ((\hat{x}_b, p_b); b \in (1, \ldots, B))$    // *Augmented labeled examples and their labels*

11: $\hat{\mathcal{U}} = ((\hat{u}_{b,k}, q_b); b \in (1, \ldots, B), k \in (1, \ldots, K))$    // *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: $\mathcal{X}' = (\text{MixUp}(\hat{\mathcal{X}}_i, \mathcal{W}_i); i \in (1, \ldots, |\hat{\mathcal{X}}|))$    // *Apply MixUp to labeled data and entries from $\mathcal{W}$*

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

15: **return** $\mathcal{X}', \mathcal{U}'$

## Implementation details

- Model = Wide ResNet-28.
- Checkpoint every $2^{16}$ training samples and report the median error rate of the last 20 checkpoints.
- 5000 examples to select hyperparameters.
- Weight decay of 0.0004.
- Datasets - CIFAR-10, CIFAR-100, STL-10, and SVHN.

## Hyperparameters

- Sharpening temperature T.
- Number of unlabled augmentations K.
- $\alpha$ for the Beta distribution used in Mixup.
- Unsupervised loss weight $\lambda_{\mathcal{U}}$.
- They find in practice that most of MixMatch's hyperparameters can be fixed. Specifically, for all experiments set T = 0.5 and K = 2.
- $\alpha = 0.75, \lambda_{\mathcal{U}} = 100$ good starting point.
- Linearly ramp up $\lambda_{\mathcal{U}}$ to its maximum value over the first 16,000 steps of training.

| 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).

| Method | 1000 labels | 5000 labels |
|---|---|---|
| CutOut [12] | - | 12.74 |
| IIC [20] | - | 11.20 |
| SWWAE [48] | 25.70 | - |
| CC-GAN[2] [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.

| 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 |

| Methods/Labels | 250 | 500 | 1000 | 2000 | 4000 |
|---|---|---|---|---|---|
| PiModel | $17.65 \pm 0.27$ | $11.44 \pm 0.39$ | $8.60 \pm 0.18$ | $6.94 \pm 0.27$ | $5.57 \pm 0.14$ |
| PseudoLabel | $21.16 \pm 0.88$ | $14.35 \pm 0.37$ | $10.19 \pm 0.41$ | $7.54 \pm 0.27$ | $5.71 \pm 0.07$ |
| Mixup | $39.97 \pm 1.89$ | $29.62 \pm 1.54$ | $16.79 \pm 0.63$ | $10.47 \pm 0.48$ | $7.96 \pm 0.14$ |
| VAT | $8.41 \pm 1.01$ | $7.44 \pm 0.79$ | $5.98 \pm 0.21$ | $4.85 \pm 0.23$ | $4.20 \pm 0.15$ |
| MeanTeacher | $6.45 \pm 2.43$ | $3.82 \pm 0.17$ | $3.75 \pm 0.10$ | $3.51 \pm 0.09$ | $3.39 \pm 0.11$ |
| MixMatch | $3.78 \pm 0.26$ | $3.64 \pm 0.46$ | $3.27 \pm 0.31$ | $3.04 \pm 0.13$ | $2.89 \pm 0.06$ |

Table 6: Error rate (%) for SVHN.

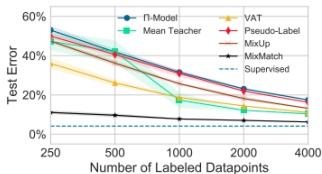| Methods/Labels | 250 | 500 | 1000 | 2000 | 4000 |
|---|---|---|---|---|---|
| PiModel | $13.71 \pm 0.32$ | $10.78 \pm 0.59$ | $8.81 \pm 0.33$ | $7.07 \pm 0.19$ | $5.70 \pm 0.13$ |
| PseudoLabel | $17.71 \pm 0.78$ | $12.58 \pm 0.59$ | $9.28 \pm 0.38$ | $7.20 \pm 0.18$ | $5.56 \pm 0.27$ |
| Mixup | $33.03 \pm 1.29$ | $24.52 \pm 0.59$ | $14.05 \pm 0.79$ | $9.06 \pm 0.55$ | $7.27 \pm 0.12$ |
| VAT | $7.44 \pm 1.38$ | $7.37 \pm 0.82$ | $6.15 \pm 0.53$ | $4.99 \pm 0.30$ | $4.27 \pm 0.30$ |
| MeanTeacher | $2.77 \pm 0.10$ | $2.75 \pm 0.07$ | $2.69 \pm 0.08$ | $2.60 \pm 0.04$ | $2.54 \pm 0.03$ |
| MixMatch | $2.22 \pm 0.08$ | $2.17 \pm 0.07$ | $2.18 \pm 0.06$ | $2.12 \pm 0.03$ | $2.07 \pm 0.05$ |

Table 7: Error rate (%) for SVHN+Extra.

Figure 2: Error rate comparison of MixMatch to baseline methods on CIFAR-10 for a varying number of labels. Exact numbers are provided in table 5 (appendix). "Supervised" refers to training with all 50000 training examples and no unlabeled data. With 250 labels MixMatch reaches an error rate comparable to next-best method's performance with 4000 labels.
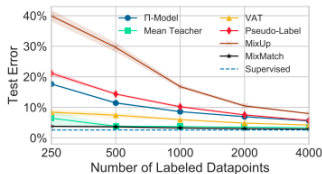
Figure 3: Error rate comparison of MixMatch to baseline methods on SVHN for a varying number of labels. Exact numbers are provided in table 6 (appendix). "Supervised" refers to training with all 73257 training examples and no unlabeled data. With 250 examples MixMatch nearly reaches the accuracy of supervised training for this model.
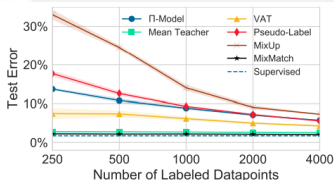
Figure 4: Error rate comparison of MixMatch to baseline methods on SVHN+Extra for a varying number of labels. With 250 examples we reach nearly the state of the art compared to supervised training for this model.

| 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.

Questions?

# Bibliography

- MixMatch: A Holistic Approach to Semi-Supervised Learning
  https://arxiv.org/pdf/1905.02249.pdf

- mixup: Beyond Empirical Risk Minimization
  https://arxiv.org/abs/1710.09412

- Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks - https://www.researchgate.net/publication/280581078$_P seudo-Label_T he_s imple_a nd_E fficient_s emi-Supervised_L earning_M ethod_f or_D eep_N eural_N etworks$