

# Semi-Supervised Learning with Meta-Gradient

Xin-Yu Zhang, Taihong Xiao, Haolin Jia, Ming-Ming Cheng,  
Ming-Hsuan Yang

*xinyuzhang@mail.nankai.edu.cn*

## 1 Prerequisites

- Consistency-based SSL

## 2 Meta-Learning Algorithm

- Overview
- Derivation
- Convergence Theory
- Trick 1 – First-Order Approximation
- Trick 2 – Mix-Up Augmentation

## 3 Experimental Results

## 4 Future Perspective

# Semi-Supervised Learning

**Semi-supervised learning (SSL):** labeled data + unlabeled data  $\Rightarrow$  better generalization ability.

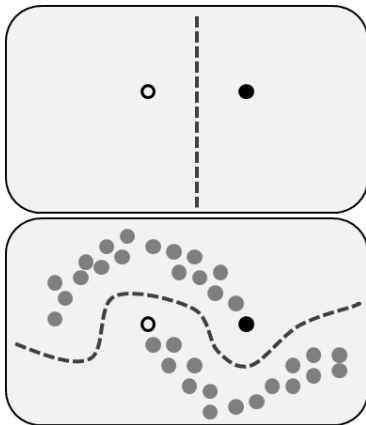


Figure: Illustration of the role of unlabeled data<sup>1</sup>.

<sup>1</sup>[https://en.wikipedia.org/wiki/Semi-supervised\\_learning](https://en.wikipedia.org/wiki/Semi-supervised_learning)

# Consistency-Based SSL

**Basic assumption:** prediction consistency against perturbations of the input signals or model weights.

---

<sup>2</sup>Miyato *et al.* Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning. In *TPAMI*, 2018

<sup>3</sup>Tarvainen & Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *NeurIPS*, 2017

# Consistency-Based SSL

**Basic assumption:** prediction consistency against perturbations of the input signals or model weights.

Two research directions for consistency-based SSL: **perturbing in the adversarial directions**<sup>2</sup> and **finding better “role models”**<sup>3</sup>.

---

<sup>2</sup>Miyato *et al.* Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning. In *TPAMI*, 2018

<sup>3</sup>Tarvainen & Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *NeurIPS*, 2017

# Consistency-Based SSL

**Basic assumption:** prediction consistency against perturbations of the input signals or model weights.

Two research directions for consistency-based SSL: **perturbing in the adversarial directions**<sup>2</sup> and **finding better “role models”**<sup>3</sup>.

**Problems:** consistency loss is rather generic and does not exploit the label information (not designed for the specific task).

---

<sup>2</sup>Miyato *et al.* Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning. In *TPAMI*, 2018

<sup>3</sup>Tarvainen & Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *NeurIPS*, 2017

# Consistency-Based SSL

**Basic assumption:** prediction consistency against perturbations of the input signals or model weights.

Two research directions for consistency-based SSL: **perturbing in the adversarial directions**<sup>2</sup> and **finding better “role models”**<sup>3</sup>.

**Problems:** consistency loss is rather generic and does not exploit the label information (not designed for the specific task).

Our work is to bridge the gap between the consistency loss and the label information.

---

<sup>2</sup>Miyato *et al.* Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning. In *TPAMI*, 2018

<sup>3</sup>Tarvainen & Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *NeurIPS*, 2017

**Challenge.** The consistency loss, which is calculated from the unlabeled data, seems to have no relationship with the labeled data.



# Overview

**Challenge.** The consistency loss, which is calculated from the unlabeled data, seems to have no relationship with the labeled data.

**Solution.** Magic of *meta-learning*: unfolding and differentiating through one SGD step.

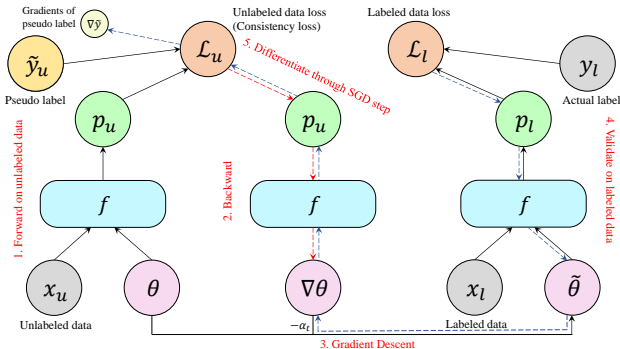


Figure: Illustration of the meta-learning philosophy.

Formulation:

$$\begin{aligned} \min_{\mathcal{Y}} \quad & \sum_{k=1}^{N^l} \mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \boldsymbol{\theta}^*(\mathcal{Y})) \\ \text{s.t. } \quad & \boldsymbol{\theta}^*(\mathcal{Y}) = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{N^u} \mathcal{L}(\mathbf{x}_i^u, \hat{\mathbf{y}}_i; \boldsymbol{\theta}). \end{aligned} \tag{1}$$

Formulation:

$$\begin{aligned} \min_{\mathcal{Y}} \quad & \sum_{k=1}^{N^l} \mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \boldsymbol{\theta}^*(\mathcal{Y})) \\ \text{s.t. } \quad & \boldsymbol{\theta}^*(\mathcal{Y}) = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{N^u} \mathcal{L}(\mathbf{x}_i^u, \hat{\mathbf{y}}_i; \boldsymbol{\theta}). \end{aligned} \tag{1}$$

Solving Eq. (1) exactly is impossible, we adopt online approximation on the batch level.

# Derivation

Initialize pseudo-labels:

$$\tilde{\mathbf{y}}_i = f(\mathbf{x}_i^u; \boldsymbol{\theta}_t). \quad (2)$$

# Derivation

Initialize pseudo-labels:

$$\tilde{\mathbf{y}}_i = f(\mathbf{x}_i^u; \boldsymbol{\theta}_t). \quad (2)$$

Compute the unlabeled data loss and back-propagate the gradients:

$$\begin{aligned} \mathcal{L}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t) &= \Phi(f(\mathbf{x}_i^u; \boldsymbol{\theta}_t), \tilde{\mathbf{y}}_i), \\ \nabla \boldsymbol{\theta}_t &= \frac{1}{B^u} \sum_{i=1}^{B^u} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t). \end{aligned} \quad (3)$$

# Derivation

Initialize pseudo-labels:

$$\tilde{\mathbf{y}}_i = f(\mathbf{x}_i^u; \boldsymbol{\theta}_t). \quad (2)$$

Compute the unlabeled data loss and back-propagate the gradients:

$$\begin{aligned} \mathcal{L}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t) &= \Phi(f(\mathbf{x}_i^u; \boldsymbol{\theta}_t), \tilde{\mathbf{y}}_i), \\ \nabla \boldsymbol{\theta}_t &= \frac{1}{B^u} \sum_{i=1}^{B^u} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t). \end{aligned} \quad (3)$$

Apply one SGD step on the model parameters:

$$\tilde{\boldsymbol{\theta}}_{t+1} = \boldsymbol{\theta}_t - \alpha_t \nabla \boldsymbol{\theta}_t, \quad (4)$$

where  $\alpha_t$  is the learning rate of the inner loop.

# Derivation

Evaluate on the labeled data and differentiate the labeled data loss:

$$\begin{aligned}\mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1}) &= \Phi(f(\mathbf{x}_k^l; \tilde{\boldsymbol{\theta}}_{t+1}), \mathbf{y}_k), \\ \nabla \tilde{\mathbf{y}}_i &= \frac{1}{B^l} \sum_{k=1}^{B^l} \nabla_{\tilde{\mathbf{y}}_i} \mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1}).\end{aligned}\tag{5}$$

# Derivation

Evaluate on the labeled data and differentiate the labeled data loss:

$$\begin{aligned}\mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1}) &= \Phi(f(\mathbf{x}_k^l; \tilde{\boldsymbol{\theta}}_{t+1}), \mathbf{y}_k), \\ \nabla \tilde{\mathbf{y}}_i &= \frac{1}{B^l} \sum_{k=1}^{B^l} \nabla_{\tilde{\mathbf{y}}_i} \mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1}).\end{aligned}\tag{5}$$

Perform one SGD step on the pseudo-labels:

$$\hat{\mathbf{y}}_i = \tilde{\mathbf{y}}_i - \beta_t \nabla \tilde{\mathbf{y}}_i,\tag{6}$$

where  $\beta_t$  is the meta learning rate,



# Derivation

Evaluate on the labeled data and differentiate the labeled data loss:

$$\begin{aligned}\mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1}) &= \Phi(f(\mathbf{x}_k^l; \tilde{\boldsymbol{\theta}}_{t+1}), \mathbf{y}_k), \\ \nabla \tilde{\mathbf{y}}_i &= \frac{1}{B^l} \sum_{k=1}^{B^l} \nabla_{\tilde{\mathbf{y}}_i} \mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1}).\end{aligned}\tag{5}$$

Perform one SGD step on the pseudo-labels:

$$\hat{\mathbf{y}}_i = \tilde{\mathbf{y}}_i - \beta_t \nabla \tilde{\mathbf{y}}_i,\tag{6}$$

where  $\beta_t$  is the meta learning rate,

Compute the consistency loss from the unlabeled data and the updated pseudo-labels.

# Meta-Learning Algorithm

---

**Algorithm 1** Meta-Learning Algorithm.

---

**Input:** regular learning rates  $\{\alpha_t\}$ ,

meta learning rates  $\{\beta_t\}$

**for**  $t := 1$  to  $\#iters$  **do**

$$\{(\mathbf{x}_k^l, \mathbf{y}_k)\}_{k=1}^{B^l} \leftarrow \text{BatchSampler}(\mathcal{D}^l)$$

$$\{\mathbf{x}_i^u\}_{i=1}^{B^u} \leftarrow \text{BatchSampler}(\mathcal{D}^u)$$

$$\tilde{\mathbf{y}}_i = f(\mathbf{x}_i^u; \boldsymbol{\theta}_t)$$

$$\mathcal{L}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t) = \Phi(f(\mathbf{x}_i^u; \boldsymbol{\theta}_t), \tilde{\mathbf{y}}_i)$$

$$\nabla \boldsymbol{\theta}_t = \frac{1}{B^u} \sum_{i=1}^{B^u} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t)$$

$$\tilde{\boldsymbol{\theta}}_{t+1} = \boldsymbol{\theta}_t - \alpha_t \nabla \boldsymbol{\theta}_t$$

$$\mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1}) = \Phi(f(\mathbf{x}_k^l; \tilde{\boldsymbol{\theta}}_{t+1}), \mathbf{y}_k)$$

$$\nabla \tilde{\mathbf{y}}_i = \frac{1}{B^l} \sum_{k=1}^{B^l} \nabla_{\tilde{\mathbf{y}}_i} \mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1})$$

$$\hat{\mathbf{y}}_i = \tilde{\mathbf{y}}_i - \beta_t \nabla \tilde{\mathbf{y}}_i$$

$$\mathcal{L}(\mathbf{x}_i^u, \hat{\mathbf{y}}_i; \boldsymbol{\theta}_t) = \Phi(f(\mathbf{x}_i^u; \boldsymbol{\theta}_t), \hat{\mathbf{y}}_i)$$

$$\nabla \hat{\boldsymbol{\theta}}_t = \frac{1}{B^u} \sum_{i=1}^{B^u} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}_i^u, \hat{\mathbf{y}}_i; \boldsymbol{\theta}_t)$$

$$\boldsymbol{\theta}_{t+1} = \text{Optimizer}(\boldsymbol{\theta}_t, \nabla \hat{\boldsymbol{\theta}}_t, \alpha_t)$$

**end**

# Convergence Theory — Convergence Guarantee

## Theorem

Let

$$G(\boldsymbol{\theta}; \mathcal{D}^l) = \frac{1}{N^l} \sum_{k=1}^{N^l} \mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \boldsymbol{\theta}_t) \quad (7)$$

be the loss function of the labeled examples. Assume

- (i)  $\nabla_{\boldsymbol{\theta}} G$  is Lipschitz-continuous with a Lipschitz constant  $L_0$ ; and
- (ii) the norm of the Jacobian matrix of  $f$  w.r.t.  $\boldsymbol{\theta}$  is upper-bounded, i.e.,

$$\exists M \in \mathbb{R}, \text{ s.t. } \|J_{\boldsymbol{\theta}} f(\mathbf{x}_i^u; \boldsymbol{\theta})\| \leq M, \quad \forall i \in \{1, \dots, N^u\}. \quad (8)$$

If the regular learning rate  $\alpha_t$  and meta learning rate  $\beta_t$  satisfy  $\alpha_t^2 \beta_t < (4M^2 L_0)^{-1}$ , then each SGD step of Alg. 1 will lead to the decrease of the validation loss  $G(\boldsymbol{\theta})$ , regardless of the selected unlabeled examples, i.e.,

$$G(\boldsymbol{\theta}_{t+1}) \leq G(\boldsymbol{\theta}_t), \quad \text{for each } t. \quad (9)$$

## Theorem

Assume the same conditions as in Thm. 10, and

$$\inf_t (\beta_t - 4\alpha_t^2\beta_t^2M^2L_0) = D_1 > 0 \quad \text{and} \quad \inf_t \alpha_t = D_2 > 0. \quad (10)$$

We further assume that the unlabeled dataset contains the labeled dataset, *i.e.*,  $\mathcal{D}^l \subseteq \mathcal{D}^u$ . Then, Alg. 1 achieves  $\mathbb{E} [\|\nabla_{\theta} G(\theta_t)\|^2] \leq \epsilon$  in  $O(1/\epsilon^2)$  steps, *i.e.*,

$$\min_{1 \leq t \leq T} \mathbb{E} [\|\nabla_{\theta} G(\theta_t)\|^2] \leq \frac{C}{\sqrt{T}}, \quad (11)$$

where  $C$  is a constant independent of the training process.

# Trick 1 – First-Order Approximation

Differentiation through the SGD step in Eq. (5) is computationally expensive since the second-order derivative is involved. By applying the chain rule, we have

$$\frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1}) = -\frac{\alpha_t}{B^u} \nabla_{\boldsymbol{\theta}}^{\top} \frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t) \cdot \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \boldsymbol{\theta}_t). \quad (12)$$

# Trick 1 – First-Order Approximation

Differentiation through the SGD step in Eq. (5) is computationally expensive since the second-order derivative is involved. By applying the chain rule, we have

$$\frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_k^l, \mathbf{y}_k; \tilde{\boldsymbol{\theta}}_{t+1}) = -\frac{\alpha_t}{B^u} \nabla_{\boldsymbol{\theta}}^{\top} \frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t) \cdot \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}_k^l, \mathbf{y}_k; \boldsymbol{\theta}_t). \quad (12)$$

Therefore, the gradients of the pseudo-labels can be formulated as

$$\nabla \tilde{y}_{i,j} = \frac{1}{B^l} \sum_{k=1}^{B^l} \frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_k^l, \mathbf{y}_k; \boldsymbol{\theta}_t) = -\frac{\alpha_t}{B^u} \nabla_{\boldsymbol{\theta}}^{\top} \frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t) \cdot \nabla \boldsymbol{\theta}_t^l, \quad (13)$$

where

$$\nabla \boldsymbol{\theta}_t^l = \frac{1}{B^l} \sum_{i=1}^{B^l} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}_i^l, \mathbf{y}_i; \boldsymbol{\theta}_t), \quad (14)$$

# Trick 1 – First-Order Approximation

By Taylor expansion,

$$\nabla \tilde{y}_{i,j} \approx -\frac{\alpha_t}{2B^u \epsilon} \left( \frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t + \epsilon \nabla \boldsymbol{\theta}_t^l) - \frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t - \epsilon \nabla \boldsymbol{\theta}_t^l) \right), \quad (15)$$

where  $\epsilon$  is a sufficiently small number.

# Trick 1 – First-Order Approximation

By Taylor expansion,

$$\nabla \tilde{y}_{i,j} \approx -\frac{\alpha_t}{2B^u\epsilon} \left( \frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t + \epsilon \nabla \boldsymbol{\theta}_t^l) - \frac{\partial \mathcal{L}}{\partial \tilde{y}_{i,j}}(\mathbf{x}_i^u, \tilde{\mathbf{y}}_i; \boldsymbol{\theta}_t - \epsilon \nabla \boldsymbol{\theta}_t^l) \right), \quad (15)$$

where  $\epsilon$  is a sufficiently small number.

Furthermore, for the MSE loss  $\Phi^{\text{MSE}}(\mathbf{p}, \mathbf{y}) = \|\mathbf{p} - \mathbf{y}\|_2^2$ , the gradients are given by:

$$\nabla \tilde{\mathbf{y}}_i \approx \frac{\alpha_t}{B^u\epsilon} \left( f(\mathbf{x}_i^u; \boldsymbol{\theta}_t + \epsilon \nabla \boldsymbol{\theta}_t^l) - f(\mathbf{x}_i^u; \boldsymbol{\theta}_t - \epsilon \nabla \boldsymbol{\theta}_t^l) \right). \quad (16)$$

(Please refer to our paper for the technical details.)



## Trick 2 – Mix-Up Augmentation

Due to the scarcity of labeled data, Mix-Up is a commonly-used technique to augment the training data and regularize the learned model.

## Trick 2 – Mix-Up Augmentation

Due to the scarcity of labeled data, Mix-Up is a commonly-used technique to augment the training data and regularize the learned model.

Assuming the labeled data and unlabeled data are of equal batch size, *i.e.*,  $B^l = B^u = B$ , then we adopt the cross-domain mixup, namely,

$$\begin{aligned} \mathbf{x}_i^{\text{in}} &= \lambda_i \mathbf{x}_i^l + (1 - \lambda_i) \mathbf{x}_i^u, \\ \mathbf{y}_i^{\text{in}} &= \lambda_i \mathbf{y}_i + (1 - \lambda_i) \hat{\mathbf{y}}_i, \end{aligned} \quad i = 1, 2, \dots, B, \quad (17)$$

where  $\lambda_1, \lambda_2, \dots, \lambda_B \stackrel{\text{i.i.d.}}{\sim} \text{Beta}(\gamma, \gamma)$ .

## Trick 2 – Mix-Up Augmentation

Due to the scarcity of labeled data, Mix-Up is a commonly-used technique to augment the training data and regularize the learned model.

Assuming the labeled data and unlabeled data are of equal batch size, *i.e.*,  $B^l = B^u = B$ , then we adopt the cross-domain mixup, namely,

$$\begin{aligned} \mathbf{x}_i^{\text{in}} &= \lambda_i \mathbf{x}_i^l + (1 - \lambda_i) \mathbf{x}_i^u, \\ \mathbf{y}_i^{\text{in}} &= \lambda_i \mathbf{y}_i + (1 - \lambda_i) \hat{\mathbf{y}}_i, \end{aligned} \quad i = 1, 2, \dots, B, \quad (17)$$

where  $\lambda_1, \lambda_2, \dots, \lambda_B \stackrel{\text{i.i.d.}}{\sim} \text{Beta}(\gamma, \gamma)$ .

Note that  $\hat{\mathbf{y}}_i$  here must be the *updated* pseudo label in Eq. (6), as the label information is encoded therein.

# Overall Algorithm

---

**Algorithm 2** Algorithm with Mix-Up Augmentation.

---

**Input:** regular learning rates  $\{\alpha_t\}$ ,

meta learning rates  $\{\beta_t\}$

**for**  $t := 1$  to  $\#iters$  **do**

$$\{(\mathbf{x}_k^l, \mathbf{y}_k)\}_{k=1}^B \leftarrow \text{BatchSampler}(\mathcal{D}^l)$$

$$\{\mathbf{x}_i^u\}_{i=1}^B \leftarrow \text{BatchSampler}(\mathcal{D}^u)$$

$$\tilde{\mathbf{y}}_i = f(\mathbf{x}_i^u; \boldsymbol{\theta}_t)$$

$$\mathcal{L}^{\text{KL}}(\mathbf{x}_k^l, \mathbf{y}_k; \boldsymbol{\theta}_t) = \Phi^{\text{KL}}(f(\mathbf{x}_k^l; \boldsymbol{\theta}_t), \mathbf{y}_k)$$

$$\nabla \boldsymbol{\theta}_t^l = \frac{1}{B} \sum_{k=1}^B \nabla_{\boldsymbol{\theta}} \mathcal{L}^{\text{KL}}(\mathbf{x}_k^l, \mathbf{y}_k; \boldsymbol{\theta}_{t+1})$$

$$\epsilon = 0.01 \|\nabla \boldsymbol{\theta}_t^l\|^{-1}$$

$$\nabla \tilde{\mathbf{y}}_i = \epsilon^{-1} (f(\mathbf{x}_i^u; \boldsymbol{\theta}_t + \epsilon \nabla \boldsymbol{\theta}_t^l) - f(\mathbf{x}_i^u; \boldsymbol{\theta}_t - \epsilon \nabla \boldsymbol{\theta}_t^l))$$

$$\hat{\mathbf{y}}_i = \tilde{\mathbf{y}}_i - \beta_t \nabla \tilde{\mathbf{y}}_i$$

$$\lambda_i \leftarrow \text{Beta}(\gamma, \gamma)$$

$$\mathbf{x}_i^{\text{in}} = \lambda \mathbf{x}_i^l + (1 - \lambda) \mathbf{x}_i^u$$

$$\mathbf{y}_i^{\text{in}} = \lambda \mathbf{y}_i + (1 - \lambda) \hat{\mathbf{y}}_i$$

$$\mathcal{L}_{\text{cls}}^{\text{KL}}(\mathbf{x}_i^{\text{in}}, \mathbf{y}_i^{\text{in}}; \boldsymbol{\theta}_t) = \Phi^{\text{KL}}(f(\mathbf{x}_i^{\text{in}}; \boldsymbol{\theta}_t), \mathbf{y}_i^{\text{in}})$$

$$\mathcal{L}_{\text{cons}}^{\text{MSE}}(\mathbf{x}_i^u, \hat{\mathbf{y}}_i; \boldsymbol{\theta}_t) = \Phi^{\text{MSE}}(f(\mathbf{x}_i^u; \boldsymbol{\theta}_t), \hat{\mathbf{y}}_i)$$

$$\nabla \hat{\boldsymbol{\theta}}_t = \frac{1}{B} \sum_{i=1}^B (\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{cls}}^{\text{KL}} + \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{cons}}^{\text{MSE}})$$

$$\boldsymbol{\theta}_{t+1} = \text{Optimizer}(\boldsymbol{\theta}_t, \nabla \hat{\boldsymbol{\theta}}_t, \alpha_t)$$

**end**

# Experiments on Small Datasets

Experiments on **SVHN**, **CIFAR-10**, and **CIFAR-100**.

Method	SVHN	CIFAR-10	CIFAR-100
PI-Model [7]	4.82%	12.36%	39.19%
TE [7]	4.42%	12.16%	38.65%
MT [9]	3.95%	12.31%	-
MT+SNTG [26]	3.86%	10.93%	-
VAT [8]	5.42%	11.36%	-
VAT+Ent [8]	3.86%	10.55%	-
VAT+Ent+SNTG [26]	3.83%	9.89%	-
VAT+VAAdD [27]	3.55%	9.22%	-
MA-DNN [28]	4.21%	11.91%	34.51%
Co-training [29]	3.29%	8.35%	34.63%
MT+fastSWA [10]	-	9.05%	33.62%
TNAR-VAE [11]	3.74%	8.85%	-
ADA-Net [24]	4.62%	10.30%	-
ADA-Net+fastSWA [24]	-	8.72%	-
DualStudent [30]	-	8.89%	32.77%
Ours	<b>3.15%</b>	<b>7.78%</b>	<b>30.74%</b>
Fully-Supervised	2.67%	4.88%	22.10%

Figure: Semi-supervised classification results.

# Experiments on ImageNet

Method	Top-1	Top-5
Labeled-Only	53.65%	31.01%
MT [9]	49.07%	23.59%
Co-training [29]	46.50%	22.73%
ADA-Net [24]	44.91%	21.18%
Ours	<b>44.87%</b>	<b>18.88%</b>
Fully-Supervised	29.15%	10.12%

Figure: Semi-supervised classification results on ImageNet.

# Experiments on ImageNet

Method	Top-1	Top-5
Labeled-Only	53.65%	31.01%
MT [9]	49.07%	23.59%
Co-training [29]	46.50%	22.73%
ADA-Net [24]	44.91%	21.18%
Ours	<b>44.87%</b>	<b>18.88%</b>
Fully-Supervised	29.15%	10.12%

Figure: Semi-supervised classification results on ImageNet.

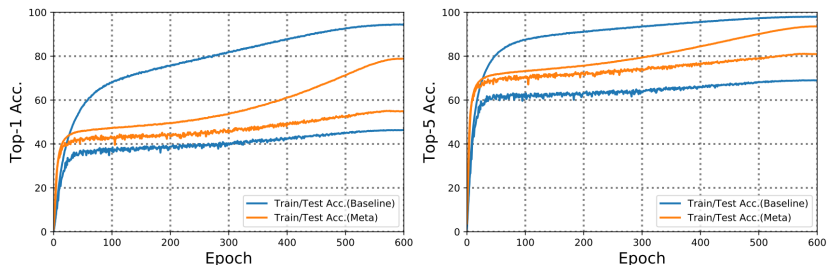


Figure: Accuracy curves of the baseline method and the meta-learning algorithm.

# Ablation Studies and Visualization

No.	Meta	Mix-Up	SVHN	CIFAR-10	CIFAR-100
1			9.76%	15.43%	38.74%
2	✓		3.68%	11.63%	35.40%
3		✓	5.60%	11.10%	32.67%
4	✓	✓	<b>3.15%</b>	<b>7.78%</b>	<b>30.74%</b>

Figure: Ablation of each component.



# Ablation Studies and Visualization

No.	Meta	Mix-Up	SVHN	CIFAR-10	CIFAR-100
1			9.76%	15.43%	38.74%
2	✓		3.68%	11.63%	35.40%
3		✓	5.60%	11.10%	32.67%
4	✓	✓	<b>3.15%</b>	<b>7.78%</b>	<b>30.74%</b>

Figure: Ablation of each component.

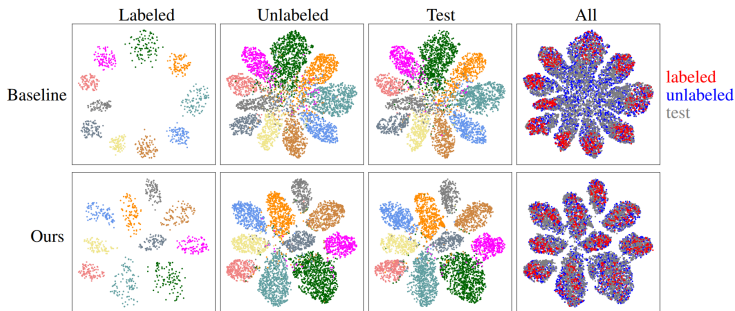


Figure: Feature Visualization.

- (i) Incorporation with the concurrent self-supervised SSL methods, e.g., FixMatch<sup>4</sup>;
- (ii) Extension beyond semi-supervised classification, e.g., weakly-supervised segmentation;
- (iii) (May be too aggressive) Generalize the current bi-level optimization into multi-level optimization.

---

<sup>4</sup>Sohn *et al.* FixMatch: Simplifying Semi-Supervised Learning with Consistency and Confidence. In: *NeurIPS*, 2020.

# Thanks!