

# A survey on deep semi-supervised learning algorithms

Ani Vanyan, Hrant Khachatrian

September 13, 2020

# Overview

- 1 Semi-supervised learning
- 2 Consistency regularization
- 3 Evaluation challenges
- 4 Multi-stage algorithms
- 5 Current state-of-the-art: Back to Basics

- **Semi-supervised learning** is about making use of unlabeled data
  - It is useful when the available labeled data is small
- Unlabeled data is cheap
  - Although not every type of unlabeled data is useful

# Consistency regularization

The idea is to make sure the neural network produces similar results for the augmented versions of the same unlabeled image.

- Definitions

- $X$  is the labeled dataset with samples  $(x, p) \in X$
- $U$  is the dataset without labels
- $f_\theta(x)$  is a function (neural network): outputs probabilities of the labels
- $Augment(x)$  is a *stochastic* function that augments the sample  $x$
- $H(\cdot, \cdot)$  is the cross entropy loss

- Consistency regularization:

- $$L_U = \frac{1}{|U|} \sum_{x \in U} \|f_\theta(Augment(x)) - f_\theta(x)\|_2^2$$

- $\Pi$ -model
  - $L = L_X + \lambda(t)L_U$
  - $L_X = \frac{1}{|X|} \sum_{(x,p) \in X} H(p, f_\theta(x))$
  - $L_U = \frac{1}{|U|} \sum_{x \in U} \|f_\theta(\text{Augment}(x)) - f_\theta(x)\|_2^2$
  - $\lambda(t)$  slowly grows from 0 to its final value  $\lambda$ 
    - which is a hyperparameter
- $\text{Augment}(x)$ 
  - Translation by  $a \sim \text{Uniform}(-2, 2)$  pixels
  - Horizontal flip (for CIFAR-10/100 only)

---

<sup>1</sup>Samuli Laine and Timo Aila. “Temporal ensembling for semi-supervised learning”. In: *ICLR* (2017).

# The problem of the unstable target

The problem with consistency loss is that it is not stable. In  $\Pi$ -model the partial solution was given. The trick is called temporal ensembling.

$$f_{temp.ens}(x) = \alpha f_{temp.ens}(Augment(x)) + (1 - \alpha) f_{\theta}(Augment(x))$$
$$L_U = \frac{1}{|U|} \sum_{x \in U} \|f_{\theta}(Augment(x)) - f_{temp.ens}(x)\|_2^2$$

The first formula is computed once per epoch.

# Mean Teacher<sup>2</sup>

- Two separate models: a Student network with  $\theta$  parameters and a Teacher with  $\theta'$  parameters.
- $L_X = \frac{1}{|X|} \sum_{(x,p) \in X} H(p, f_{\theta}^{student}(x))$
- Teacher's parameters are updated at each iteration
- $\theta'_t = \alpha \theta'_{t-1} + (1 - \alpha) \theta_t$
- $L_U = \frac{1}{|U|} \sum_{x \in U} \|f_{\theta}^{student}(Augment(x)) - f_{\theta'}^{teacher}(Augment(x))\|_2^2$
- Teacher is *not* trained via backpropagation.

---

<sup>2</sup>Antti Tarvainen and Harri Valpola. "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results". In: *Advances in neural information processing systems*. 2017, pp. 1195–1204.

# Virtual Adversarial Training<sup>3</sup>

- Regular adversarial training
  - For each  $(x, p) \in X$
  - $L_{adv} = H(p, f_{\theta}(x + r_{adv}))$
  - $r_{adv} = \arg \max_{r: \|r\| < \epsilon} H(p, f_{\theta}(x + r))$
- Virtual adversarial training
  - For each  $x \in U$
  - $L_{adv} = H(f_{\theta}(x), f_{\theta}(x + r_{adv}))$
  - $r_{adv} = \arg \max_{r: \|r\| < \epsilon} H(f_{\theta}(x), f_{\theta}(x + r))$
  - Fast approximation method for  $r_{adv}$
- Entropy minimization
  - $L_{ent} = \frac{1}{|X|+|U|} \sum_{x \in X \cup U} H(f_{\theta}(x))$

---

<sup>3</sup>Takeru Miyato et al. “Virtual adversarial training: a regularization method for supervised and semi-supervised learning”. In: *IEEE transactions on pattern analysis and machine intelligence* 41.8 (2018), pp. 1979–1993.

- Reproducing results
  - Reimplemented many known methods in a single code repository
  - WideResNet-28-2, batch normalization, leaky ReLU
  - Adam optimizer with fixed  $\beta_1$  and  $\beta_2$
  - Fixed data augmentation and preprocessing
    - Although, different for CIFAR-10 and SVHN!
  - Equal hyperparameter tuning budget!
    - 1000 trials of Gaussian-Process-based black box optim. in Google Cloud
    - Model selection on the full validation set
    - Different initial learning rates for different methods!
    - VAT's  $\epsilon$  was different on CIFAR-10 and SVHN

---

<sup>4</sup>Avital Oliver et al. “Realistic evaluation of deep semi-supervised learning algorithms”. In: *Advances in Neural Information Processing Systems*. 2018, pp. 3235–3246.

- Transfer learning
  - Transfer from (resized) ImageNet is a strong baseline, ignored in most papers
- Class distribution mismatch
  - Adding unlabeled data from a mismatched set of classes can actually **hurt** the performance compared to not using unlabeled data at all

---

<sup>5</sup>Avital Oliver et al. “Realistic evaluation of deep semi-supervised learning algorithms”. In: *Advances in Neural Information Processing Systems*. 2018, pp. 3235–3246.

- Similar underlying model, but with different  $X'$  and  $U'$ 
  - $X', U' = \text{MixMatch}(X, U, T, K, \alpha)$
  - $L_X = \frac{1}{|X'|} \sum_{(x', p') \in X'} H(p', f_\theta(x'))$
  - $L_U = \frac{1}{|U'|} \sum_{(x', q') \in U'} \|q' - f_\theta(x')\|_2^2$
  - $L = L_X + \lambda(t)L_U$

---

<sup>6</sup>David Berthelot et al. "Mixmatch: A holistic approach to semi-supervised learning".  
In: *accepted at NeurIPS'19 (2019)*.

MixMatch uses trick for obtaining new unlabeled examples called MixUp.

- Constructs new pairs by convex combination of existing pairs
- For each pair  $(x_1, p_1)$  and  $(x_2, p_2)$ :
  - Sample  $\lambda \sim \text{Beta}(\alpha, \alpha)$
  - $\lambda' = \max(\lambda, 1 - \lambda)$  to make sure it's close to 0
    - This step did not exist in the original MixUp paper
  - $x' = \lambda' x_1 + (1 - \lambda') x_2$
  - $p' = \lambda' p_1 + (1 - \lambda') p_2$
  - Return  $(x', p')$
- $\text{MixUp}(A, B)$  is a set of samples “closer” to  $A$

---

<sup>7</sup>Hongyi Zhang et al. “mixup: Beyond empirical risk minimization”. In: *arXiv preprint arXiv:1710.09412* (2017).

It is a more complicated version of MixMatch.

The two main additions are:

- Distribution alignment. The predicted probabilities on a batch of unlabeled examples are scaled to match the distribution of the labels present in the labeled subset.
- Anchored augmentation. The target label for unlabeled examples is determined using weakly augmented versions of the images, while the prediction for the same images is computed by using strongly augmented versions.

---

<sup>8</sup>David Berthelot et al. “ReMixMatch: Semi-Supervised Learning with Distribution Alignment and Augmentation Anchoring”. In: *arXiv preprint arXiv:1911.09785 / accepted at ICLR'2020* (2019).

UDA is quite similar to Virtual Adversarial Training, but replaces virtual adversarial example generation with a very strong augmentation called RandAugment<sup>9</sup>.

UDA uses a training technique called Training Signal Annealing to reduce overfitting when there is a huge gap between the amount of unlabeled data and labeled data.

---

<sup>9</sup>Ekin D Cubuk et al. "RandAugment: Practical data augmentation with no separate search". In: *arXiv preprint arXiv:1909.13719* (2019).

<sup>10</sup>Qizhe Xie et al. "Unsupervised data augmentation". In: *arXiv preprint arXiv:1904.12848* (2019).

FixMatch uses CutOut<sup>11</sup> along with RandAugment or CTAugment<sup>12</sup> as a strong augmentation procedure known from previous papers (UDA and ReMixMatch, respectively).

In contrast to UDA and other methods, FixMatch performs  $\arg \max$  on the guessed label, so it essentially becomes equivalent to pseudo-labeling. Additionally, FixMatch ignores the guessed labels if the confidence is lower than  $\tau = 0.95$  threshold.

---

<sup>11</sup>Terrance DeVries and Graham W Taylor. “Improved regularization of convolutional neural networks with cutout”. In: *arXiv preprint arXiv:1708.04552* (2017).

<sup>12</sup>David Berthelot et al. “ReMixMatch: Semi-Supervised Learning with Distribution Alignment and Augmentation Anchoring”. In: *arXiv preprint arXiv:1911.09785 / accepted at ICLR'2020* (2019).

<sup>13</sup>Kihyuk Sohn et al. “Fixmatch: Simplifying semi-supervised learning with consistency and confidence”. In: *arXiv preprint arXiv:2001.07685* (2020).

# Results

Method	Model	Parameters	Error (%)	
			CIFAR-10(4k)	SVHN(1k)
$\Pi$ -Model (2016)	Conv-Large	3.1M	$12.36 \pm 0.31$	$4.82 \pm 0.17$
Mean Teacher (2017)	Conv-Large	3.1M	$12.31 \pm 0.28$	$3.95 \pm 0.19$
VAT + EntMin (2017)	Conv-Large	3.1M	$10.55 \pm 0.05$	$3.86 \pm 0.11$
Mean Teacher (2017)	Shake-Shake	26M	$6.28 \pm 0.15$	-
UDA(RandAugment) (2019)	Shake-Shake	26M	3.7	-
MixMatch (2019)	WRN	26M	$4.95 \pm 0.08$	-
UDA(RandAugment) (2019)	WRN-28-2	1.5M	$5.29 \pm 0.25$	<b><math>2.55 \pm 0.09</math></b>
ReMixMatch (2019)	WRN-28-2	1.5M	$5.14 \pm 0.04$	$2.83 \pm 0.30$
FixMatch(RA) (2020)	WRN-28-2	1.5M	<b><math>4.26 \pm 0.05</math></b>	<b><math>2.28 \pm 0.11</math></b>
FixMatch(CTA) (2020)	WRN-28-2	1.5M	<b><math>4.31 \pm 0.15</math></b>	<b><math>2.36 \pm 0.19</math></b>
UDA(RandAugment) (2019)	PyramidNet	26M	<b>2.7</b>	-

Thanks.