

MixMatch: A Holistic Approach to Semi-Supervised Learning

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

Google Research

NIPS 2019, 193회 인용

Park, MinKyu

2020.07.08

Dongguk University

Artificial Intelligence Laboratory

mkpark73@dongguk.edu

Variations on MixMatch

• These include:

1) ReMixMatch: Semi-Supervised Learning with Distribution Alignment and Augmentation Anchoring <u>https://arxiv.org/abs/1911.09785</u>

2) MetaMixUp: Learning Adaptive Interpolation Policy of MixUp with Meta-Learning <u>https://arxiv.org/abs/1908.10059</u>

3) EnAET: Self-Trained Ensemble AutoEncoding Transformations for Semi-Supervised Learning <u>https://arxiv.org/abs/1911.09265</u>

4) RealMix: Towards Realistic Semi-Supervised Deep Learning Algorithms <u>https://arxiv.org/abs/1912.08766</u>

- Key prerequisite reference
 - mixup: Beyond Empirical Risk Minimization <u>https://arxiv.org/abs/1710.09412</u>

Semi-Supervised Learning

Using labelled as well as unlabelled data to perform certain learning tasks



Figure 3. Illustration of the usefulness of unlabeled data.

SSL: Representative approaches

Generative methods

 Using a generative model for the classifier and employing EM to model the label estimation or parameter estimation process

S3VMs (Semi-Supervised SVMs)

- Using unlabeled data to adjust the decision boundary such that it goes through the less dense region

Graph-based methods

Using unlabeled data to regularize the learning process via graph regularization

Disagreement-based methods

 multiple learners are trained for the task and the disagreements among the learners are exploited during the SSL process

Semi-supervised clustering

 a technique that partitions unlabeled data by making use of domain knowledge, usually expressed as pairwise constraints among instances or just as an additional set of labeled instances

Semi-Supervised Learning by Disagreement

Co-training



Fig. 1. An illustration of the co-training procedure

Entropy Minimization

encourages the model to output confident predictions on unlabeled data



a generic model $p_{model}(y \mid x; \theta)$ which produces a distribution over class labels y for an input x with parameters θ .

- Entropy Minimization
 - encourages the model to output confident predictions on unlabeled data





Low entropy for low temperature $(T \rightarrow 0)$



Consistency Regularization

 encourages the model to produce the same output distribution when its inputs are perturbed



Generic Regularization (Traditional Regularization)

cat dog

- encourages the model to generalize well and avoid overfitting the training data.



[0.7, 0.3] cat dog

cat dog

- Generic Regularization (Traditional Regularization)
 - use weight decay which penalizes the L2 norm of the model parameters
 - MixUp : Encourage the model to have strictly linear behavior between examples.

 $\tilde{x} = \lambda x_i + (1 - \lambda) x_j$, where x_i, x_j are raw input vectors $\tilde{y} = \lambda y_i + (1 - \lambda) y_j$, where y_i, y_j are one-hot label encodings

```
# y1, y2 should be one-hot vectors
for (x1, y1), (x2, y2) in zip(loader1, loader2):
    lam = numpy.random.beta(alpha, alpha)
    x = Variable(lam * x1 + (1. - lam) * x2)
    y = Variable(lam * y1 + (1. - lam) * y2)
    optimizer.zero_grad()
    loss(net(x), y).backward()
    optimizer.step()
```

(a) One epoch of mixup training in PyTorch.



(b) Effect of *mixup* ($\alpha = 1$) on a toy problem. Green: Class 0. Orange: Class 1. Blue shading indicates p(y = 1|x).

MixMatch: How to combine three concepts?



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.

$$\mathcal{X}', \mathcal{U}' = \operatorname{MixMatch}(\mathcal{X}, \mathcal{U}, T, K, \alpha)$$
 (2)

$$\mathcal{L}_{\mathcal{X}} = \frac{1}{|\mathcal{X}'|} \sum_{x, p \in \mathcal{X}'} \mathrm{H}(p, \mathrm{p}_{\mathrm{model}}(y \mid x; \theta))$$
(3)

$$\mathcal{L}_{\mathcal{U}} = \frac{1}{L|\mathcal{U}'|} \sum_{u,q \in \mathcal{U}'} \|q - p_{\text{model}}(y \mid u; \theta)\|_2^2$$
(4)

$$\mathcal{L} = \mathcal{L}_{\mathcal{X}} + \lambda_{\mathcal{U}} \mathcal{L}_{\mathcal{U}}$$
(5)

where H(p,q) is the cross-entropy between distributions p and q, and T, K, α , and $\lambda_{\mathcal{U}}$ are hyperparameters described below. The full MixMatch algorithm is provided in algorithm 1, and a diagram of the label guessing process is shown in fig. 1. Next, we describe each part 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 $\mathcal{X} = ((x_b, p_b); b \in (1, ..., B))$, batch of unlabeled examples $\mathcal{U} = (u_b; b \in (1, ..., B))$, sharpening temperature T, number of augmentations K, Beta distribution parameter α for MixUp.

2: for b = 1 to *B* do

- 3: $\hat{x}_b = \operatorname{Augment}(x_b)$ // 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: $\bar{q}_b = \frac{1}{K} \sum_k p_{\text{model}}(y \mid \hat{u}_{b,k}; \theta) / Compute average predictions across all augmentations of <math>u_b$ 8: $q_b = \text{Sharpen}(\bar{q}_b, T) / Apply temperature sharpening to the average prediction (see eq. (7))$
- 9: end for
- 10: $\hat{\mathcal{X}} = ((\hat{x}_b, p_b); b \in (1, ..., B))$ // 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: $\mathcal{X}' = (\operatorname{MixUp}(\hat{\mathcal{X}}_i, \mathcal{W}_i); i \in (1, \dots, |\hat{\mathcal{X}}|)) / Apply \operatorname{MixUp}$ to labeled data and entries from \mathcal{W}
- 14: $\mathcal{U}' = \left(\operatorname{MixUp}(\hat{\mathcal{U}}_i, \mathcal{W}_{i+|\hat{\mathcal{X}}|}); i \in (1, \dots, |\hat{\mathcal{U}}|)\right) / / Apply \operatorname{MixUp}$ to unlabeled data and the rest of \mathcal{W} 15: return $\mathcal{X}', \mathcal{U}'$

MixMatch: Data Augmentation

- Stochastic transformation of the datapoint in such a way that its label remains unchanged.
- Apply data augmentation on both labeled and unlabeled data (Rotation, Flip, etc)

- 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

MixMatch: Label Guessing

- compute the average of the model's predicted class distributions across all the K augmentations of $\mathcal{U}b$

$$\bar{q}_b = \frac{1}{K} \sum_{k=1}^{K} p_{\text{model}}(y \mid \hat{u}_{b,k}; \theta)$$

MixMatch: Sharpening

- Add Sharpening for entropy minimization
- apply a sharpening function to reduce the entropy of the label distribution
- use the common approach of adjusting the "temperature" of this categorical distribution

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

MixMatch: MixUp

• Apply MixUp both to labeled and unlabeled samples with label guessing.

$$\lambda \sim \text{Beta}(\alpha, \alpha)$$
$$\lambda' = \max(\lambda, 1 - \lambda)$$
$$x' = \lambda' x_1 + (1 - \lambda') x_2$$
$$p' = \lambda' p_1 + (1 - \lambda') p_2$$

10: $\hat{\mathcal{X}} = ((\hat{x}_b, p_b); b \in (1, ..., B)) // Augmented labeled examples and their labels$ $11: <math>\hat{\mathcal{U}} = ((\hat{u}_{b,k}, q_b); b \in (1, ..., B), k \in (1, ..., K)) // Augmented unlabeled examples, guessed labels$ $12: <math>\mathcal{W} = \text{Shuffle}(\text{Concat}(\hat{\mathcal{X}}, \hat{\mathcal{U}})) // \text{Combine and shuffle labeled and unlabeled data}$ 13: $\mathcal{X}' = (\text{MixUp}(\hat{\mathcal{X}}_i, \mathcal{W}_i); i \in (1, ..., |\hat{\mathcal{X}}|)) // \text{Apply MixUp to labeled data and entries from } \mathcal{W}$ 14: $\mathcal{U}' = (\text{MixUp}(\hat{\mathcal{U}}_i, \mathcal{W}_{i+|\hat{\mathcal{X}}|}); i \in (1, ..., |\hat{\mathcal{U}}|)) // \text{Apply MixUp to unlabeled data and the rest of } \mathcal{W}$

MixMatch: Hyper-parameters

- T : Temperature for sharpening (T = 0.5)
- K : Number of data augmentation for unlabeled samples (K = 2)
- λ_{U} : Unsupervised loss weigh (λ_{U} = 100)
- α : MixUp parameter (α = 0.75)

Experiments

- Implementation details
 - in all experiments we use the "Wide ResNet-28" model
- four standard benchmark datasets
 - CIFAR-10 and CIFAR-100 , SVHN , and STL-10
- Baseline Methods
 - П-Model (ICLR 2017),
 - Mean Teacher (NIPS 2017),
 - Virtual Adversarial Training (ICLR 2017),
 - Pseudo-Label
 - MixUp

П-Model

П-model





source : [Laine & Aila, 2017]

Mean Teacher

averaging model weights instead of predictions



* source : [Tarvainen & Harri Valpola, 2017]

Virtual Adversarial Training(VAT)



Step 1 : Generate the adversarial image



Step 2: Minimize the KL divergence



• This simple method achieves state-of-the-art performance on benchmark datasets

Ablation study: all components are important

	Ablation	250 labels	4000 labels	
	MixMatch	11.80	6.00	
Averaging	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 -	- Sharpening
Mixup	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	

References

[Laine & Aila, 2017] Temporal Ensembling for Semi-Supervised Learning, ICML2017

[Tarvainen & Harri Valpola, 2017] Mean teachers are better role models: Weight-averaged consistency targ ets improve semi-supervised deep learning results, NIPS 2017

[Miyato et al., 2017] Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning, TPAMI 2019

[Berthelot et al., 2019] MixMatch: A Holistic Approach to Semi-Supervised Learning, NIPS 2019