

COOPERATIVE LEARNING FOR NOISY SUPERVISION

Hao Wu¹, Jiangchao Yao¹, Ya Zhang^{1,2,*}, Yanfeng Wang^{1,2}

¹Cooperative Medianet Innovation Center, Shanghai Jiao Tong University, China

²Shanghai Artificial Intelligence Laboratory, China

{howiethepeanut, sunarker, ya_zhang, wangyanfeng}@sjtu.edu.cn

ABSTRACT

Learning with noisy labels has gained the enormous interest in the robust deep learning area. Recent studies have empirically disclosed that utilizing dual networks can enhance the performance of single network but without theoretic proof. In this paper, we propose Cooperative Learning (Cool) framework for noisy supervision that analytically explains the effects of leveraging dual or multiple networks. Specifically, the simple but efficient combination in Cool yields a more reliable risk minimization for unseen clean data. A range of experiments have been conducted on several benchmarks with both synthetic and real-world settings. Extensive results indicate that Cool outperforms several state-of-the-art methods.

Index Terms— deep learning, noisy supervision, cooperative learning

1. INTRODUCTION

Large-scale supervised training datasets have significantly driven the success of deep neural networks (DNNs) in computer vision [1, 2, 3]. However, accurate annotations provided by the human experts are usually expensive to collect in many real-world applications. Therefore, the alternative ways such as crowdsourcing [4] and web query [5] are explored to reduce the cost. For example, there are a plethora of images tagged by users on open platforms which can be exploited easily and inexhaustibly. The negative impact is these approaches also inevitably introduce label noise to the dataset since the user annotations are not very reliable. Considering DNNs are capable of fitting extremely noisy labels [6], it is important to make the training robust to such label noise.

Recent studies can be roughly categorized into two classes in terms of *the classifier count*, i.e., single-network structure and dual-network structure. For the former, it usually refers to robust surrogate losses, noise transition and self-paced learning. For instance, Bootstrap [7] applies the perceptual consistency to the cross-entropy loss to mitigate the influence of

label noise. Forward [8] build a transition matrix on top of the classifier to absorb the noise. Self-paced MentorNet [9] selects small loss samples as clean instances and learn only from these instances. For the latter, it introduces twin classifiers to be the teacher of each other by which boosts the classification performance of the single network. Decouple [10] utilizes the prediction disagreement of twin networks to select more informative samples as supervision. CLC [11] leverages the entropy criterion to collaboratively correct the labels. Co-distillation [12, 13] distills the knowledge of the one to supervise the other one and vice versa. Co-teaching [14] leverages two networks to select small loss instances for cross update.

Although the current dual-network structure empirically shows improvement over the single version, it lacks of theoretical analysis and guarantee that it can always work. Besides, a natural question is whether introducing more learners can further benefit the learning with noisy supervision. To explore these limits and give a more general scope, we propose a Cooperative Learning (Cool) paradigm that multiple classifiers work cooperatively for noisy supervision. Specifically, we firstly demonstrate the dual-network structure yields lower risk than that associated with the single network in some suitable cooperation. Then, we give a sufficient condition for the case of more learners, where the risk is negatively correlated to the number of classifiers. Generally, even although the classifiers are imperfect with noisy supervision, a lower risk can be achieved when the more disagreement is introduced. Finally, based on these analysis, a cooperative learning framework is introduced that the cooperation supervision is utilized to improve the performance. The main contribution can be summarized into the following three points. We demonstrate in the presence of noisy supervision, the linear combination of predictions from multiple networks yields a more reliable supervision than predictions from either single classifier in some conditions. A Cooperative Learning framework is introduced, where multiple different imperfect classifiers produce supervision cooperatively to iteratively boost the performance. We empirically verify the proposed method on CIFAR-10, CIFAR-100 with synthetic noise and three large-scale real-world datasets namely Clothing1M, Food-101N and WebVision. Comprehensive experiments show that Cool outperforms several state-of-the-art methods.

This work is supported by the National Key Research and Development Program of China (No. 2019YFB1804304), SHEITC (No. 2018-RGZN-02046), 111 plan (No. BP0719010), and STCSM (No. 18DZ2270700), and State Key Laboratory of UHD Video and Audio Production and Presentation.

*Corresponding author

2. THE PROPOSED METHOD

2.1. Preliminaries

Given a noisy dataset $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$, where N is the sample number, x_n denotes an image instance and $y_n \in \{0, 1\}^c$ is the corresponding noisy label, we target to produce more reliable supervision for the classifiers to learn in the presence of label noise. Assume we use $f_i (i = 1, 2, \dots)$ to represent a classifier with index i and p_i indicates the prediction of f_i . We explore to achieve this goal via the cooperation of multiple classifiers.

2.2. Dual-Network Cooperative Learning

As empirically indicated in several works [12, 13, 14, 10], the dual-network structure is easy to acquire a robust classification performance when learning with noisy supervision. To understand the law of this phenomenon, we deduce the theoretical analysis in light of risk minimization. We term this methodology as dual-network Cooperative Learning (CooL) to ease the explanation and unify the notion of this work.

Suppose there are two classifiers f_1 and f_2 and we use the combination of the predictions p_1 and p_2 respectively output from f_1 and f_2 as the new cooperation supervision:

$$\hat{p}^\lambda = \lambda p_1 + (1 - \lambda)p_2. \quad (1)$$

λ is the cooperation parameter to balance between the predictions from the two classifiers. To measure the reliability of a certain supervision \tilde{y} with respect to the ground-truth label y^* , we define the noisy supervision risk on the training set $r_{\tilde{y}} = E_{\mathcal{D}}[\|\tilde{y} - y^*\|^2]$.¹ In the following, we will show that with a suitable choice of λ , leveraging the cooperation supervision in Eq. (1) yields lower risk than the individual risk of the either model.

Theorem 1. *There always exists a λ that makes the risk $r_{\hat{p}^\lambda}$ of the dual-network cooperation lower than² individual risks of two non-identical networks that do not always produce the incorrect predictions at the same time, i.e.,*

$$\exists \lambda, r_{\hat{p}^\lambda} < \min\{r_{p_1}, r_{p_2}\},$$

Proof. First, the risks of the individual predictions from f_1 and f_2 are quantified respectively as the following terms,

$$r_{p_1} = E_{\mathcal{D}}[\|p_1 - y^*\|^2], r_{p_2} = E_{\mathcal{D}}[\|p_2 - y^*\|^2].$$

Then, the cooperation risk can be decomposed as follows,

$$\begin{aligned} r_{\hat{p}^\lambda} &= E_{\mathcal{D}}[\|\hat{p}^\lambda - y^*\|^2] = E_{\mathcal{D}}[\|\lambda p_1 + (1 - \lambda)p_2 - y^*\|^2] \\ &= \lambda^2 r_{p_1} + 2\lambda(1 - \lambda)r_{p_1 p_2} + (1 - \lambda)^2 r_{p_2}. \end{aligned}$$

¹We take inspiration from Li et al. [15] while our scenarios are different. We propose to measure risks on the training set with no clean labels.

²We exclude the situation where one classifier strictly dominates the other otherwise there is no need for cooperation since λ will be set to 0 or 1 and the cooperation risk will be equal to the lower one.

where $r_{p_1 p_2} \triangleq E_{\mathcal{D}_t}[(p_1 - y^*)^T(p_2 - y^*)]$. For two non-identical classifiers f_1 and f_2 that do not always produce the incorrect predictions at the same time, as p_1, p_2 are label distributions and y^* is the one-hot label, we will have

$$0 \leq r_{p_1 p_2} < \min\{r_{p_1}, r_{p_2}\}. \quad (2)$$

By setting $\nabla_{\lambda} r_{\hat{p}^\lambda} = 0$, we obtain

$$\begin{aligned} \min_{\lambda} r_{\hat{p}^\lambda} &= r_{p_1} - \frac{(r_{p_1} - r_{p_1 p_2})^2}{r_{p_1} + r_{p_2} - 2r_{p_1 p_2}} \\ &= r_{p_2} - \frac{(r_{p_2} - r_{p_1 p_2})^2}{r_{p_1} + r_{p_2} - 2r_{p_1 p_2}} \\ &< \min\{r_{p_1}, r_{p_2}\}. \end{aligned} \quad (3)$$

which concludes the proof of Theorem 1. \square

Remark 1. *From Theorem 1, we know the suitable dual-network cooperation can achieve a lower risk than the individual network. Furthermore, it can be found the minimum risk obtained in Eq. (3) is positively correlated to $r_{p_1 p_2}$. This term reflects divergences between the two classifiers on the samples where both of the classifiers make incorrect predictions. The optimal situation is that two classifiers never make mistakes at the same, then we have $r_{p_1 p_2} = 0$. Intuitively, both p_1 and p_2 are deviated from the true label y^* . However, these deviations are towards random directions in the presence of stochastic label noise, the proposed cooperation supervision can be closer to the true label.*

2.2.1. Connection and Difference.

Here, we rethink two representative dual-network methods namely Co-distillation and Co-teaching through the lens of Cooperative Learning.

Co-distillation represents a branch of studies which leverage the model predictions to rectify the noisy labels. It is the case where p_1 in Eq. (1) is substituted to the noisy label y . Correspondingly setting $\lambda = \frac{r_{p_2}}{r_y + r_{p_2}}$, the optimal risk for f_1 can be simplified as $\frac{r_y r_{p_2}}{r_y + r_{p_2}}$, which is lower than the risk r_{p_2} . This explains why Co-distillation shows improvement. Nevertheless, it also points out one defect that r_y is fixed and cannot be improved along with the decrease of r_{p_2} .

Co-teaching represents a line of studies which select samples with a certain criterion for training. Thus, we can adjust the risk by modifying the correlated dataset (removing the unreliable samples). According to [14], the supervision of f_1 is a candidate set \mathcal{D}_{c_2} selected by f_2 and vice versa. The risk for f_1 is then denoted as $E_{\mathcal{D}_{c_2}}[\|y - y^*\|^2]$. Generally, it is a lower risk than that on \mathcal{D} with the help of the small loss trick. If we linearly combine the supervision like CooL, the cooperation risk is also a linear combination with respect to λ . And the minimum will be obtained at the boundaries, i.e., the smaller one in r_{f_1} and r_{f_2} . In this case, a more reasonable way is to

choose one of the two candidate sets, which has a lower risk. In Co-teaching, they utilize both sets for cross update which may impair the performance.

2.3. Generalized Cooperative Learning

In this section, we aim to generalize our dual-network Cool to a multi-network variant, which is able to achieve a even lower risk. Given n non-identical classifiers, we denote the new cooperation supervision as $\hat{p}^\lambda = \lambda \mathbf{p}$, where λ is a n -dimension row vector with summation equal to 1 and \mathbf{p} is a stack of p_1, \dots, p_n in rows. We define $\mathbf{R} = (r_{ij})_{n \times n}$ where $r_{ij} = E_{\mathcal{D}}[(p_i - y^*)^T (p_j - y^*)]$. The corresponding diagonal elements are the individual risks associated with the predictions of the n classifiers respectively. In the following, we analyze the cooperation risk for multiple classifiers.

Theorem 2. Given the cooperation supervision $\hat{p}^\lambda = \lambda \mathbf{p}$, the associated risk is $r_{\hat{p}^\lambda} = \lambda \mathbf{R} \lambda^T$. An invertible \mathbf{R} yields,

$$\min_{\lambda} r_{\hat{p}^\lambda} = \frac{1}{\sum_{i,j=1}^n [\mathbf{R}^{-1}]_{i,j}}. \quad (4)$$

If all non-identical classifiers are independently trained in the same settings, so that the following conditions satisfy

$$\begin{aligned} r_{ii} &= r_{diag}, \forall i = 1, \dots, n \\ r_{ij} &= r_{off}, \forall i \neq j \end{aligned} \quad (5)$$

Then, the minimum cooperation risk in Eq. (4) will be

$$\min_{\lambda} r_{\hat{p}^\lambda} = \frac{1}{n} (r_{diag} - r_{off}) + r_{off} < r_{diag}, \quad (6)$$

which is lower than the individual risks of all classifiers.

Proof. The risk associated with the new supervision is

$$\begin{aligned} r_{\hat{p}^\lambda} &= E_{\mathcal{D}}[\|\lambda \mathbf{p} - y^*\|^2] = E_{\mathcal{D}}[\|\sum_{i=1}^n \lambda_i (p_i - y^*)\|^2] \\ &= \lambda \mathbf{R} \lambda^T, \text{ s.t. }, (1, \dots, 1) \lambda^T = 1. \end{aligned}$$

Leveraging the Lagrange multipliers, we now minimize,

$$\min f(\lambda, \mu) = \lambda \mathbf{R} \lambda^T - \mu ((1, \dots, 1) \lambda^T - 1)$$

By setting $\frac{\partial f(\lambda, \mu)}{\partial \lambda} = 0$ and $\frac{\partial f(\lambda, \mu)}{\partial \mu} = 0$, we obtain,

$$\lambda_0 = \frac{\mathbf{R}^{-1} (1, \dots, 1)^T}{\sum_{i,j=1}^n [\mathbf{R}^{-1}]_{i,j}}, \mu_0 = \frac{2}{\sum_{i,j=1}^n [\mathbf{R}^{-1}]_{i,j}}.$$

Thus, the minimum risk associated with $r_{\hat{p}^\lambda}$ is

$$\min_{\lambda} r_{\hat{p}^\lambda} = f(\lambda_0, \mu_0) = \frac{1}{\sum_{i,j=1}^n [\mathbf{R}^{-1}]_{i,j}}.$$

To further analyze the characteristics of above equation, if we have Eq. (5) satisfied Eq. (4) will be deduced as

$$\min_{\lambda} r_{\hat{p}^\lambda} = \frac{1}{n} (r_{diag} - r_{off}) + r_{off} < r_{diag}.$$

which concludes the proof of Theorem 2. \square

Remark 2. From Theorem 2, we can see that the first term on the RHS of Eq. (6) indicates by leveraging more classifiers (increasing n), we can monotonically obtain a lower risk. In this case, \mathbf{R} being diagonally dominant yields a necessary and sufficient condition where the risk is inversely proportional to the number of the classifiers. Besides, similar to the claim in Remark 1, the off-diagonal element r_{off} in \mathbf{R} characterizes the divergence of two networks. If two classifiers are complementary with each other, they will work better cooperatively even though they are imperfect.

2.4. The Cooperative Learning Framework

The theoretical analysis in previous section tells us that the cooperation of multiple classifiers can lower the supervision risk and the lower bound is determined by the divergence between the classifiers. Based on this, we introduce a new Cooperative Learning (Cool) framework where the proposed cooperation supervision namely the combination of the predictions from the multiple classifiers is adopted to re-train the individual networks. As claimed in Remark 1 and 2, the prerequisite of a better performance on noisy datasets via cooperation, is to generate diverse classifiers. We thus make the classifiers learn from different sources of information to construct pattern bias. Specifically, we pre-train the n classifiers respectively on $\mathcal{D}_1, \dots, \mathcal{D}_n$, the different n partitions of D . Note this pre-training style relies on the assumption that the subset is still sufficient enough to learn a classifier exhibiting the same risk on the whole dataset. Thus we are not able to infinitely add classifiers to lower the risk due to limited data.

After obtaining multiple different pre-trained classifiers, we can utilize the combination of their predictions to train better classifiers. Instead of training another student network with such supervision like [15], we iteratively train the classifiers with the objective function as follows,

$$L(f_i) = l_{\mathcal{D}}(\hat{p}^\lambda, f_i) + \alpha l_{\mathcal{D}_i}(y, f_i) + \beta h_{\mathcal{D}_i}(f_i). \quad (7)$$

In the first term of Eq. (7), the network f_i is supervised by the cooperation supervision \hat{p}^λ , which is the key module of this paper. The second part is an auxiliary term that supervises f_i with the original labels in the early phase but will be gradually canceled out as the model is capable of memorizing the noisy labels. The third term is the entropy of the model predictions which prevents the output of the network f_i from degenerating to the uniform distribution. As for the hyperparameters α and β , we empirically assign small weights like [13]. The complete training process is summarized in Algorithm 1.

Algorithm 1 The Cool Algorithm

Require: A noisy set \mathcal{D} , multiple networks $f_1, \dots, f_n, \lambda, \alpha, \beta$.
1: Randomly partition \mathcal{D} into $\mathcal{D}_1, \dots, \mathcal{D}_n$.
2: Directly pre-train f_i on \mathcal{D}_i ($i = 1, 2, \dots$) respectively.
3: **for** epoch $i = \text{StartEpoch}$ to MaxEpoch **do**
4: **for** batch $j = 1$ to $\frac{|\mathcal{D}|}{\text{BatchSize}}$ **do**
5: **for** $k = 1$ to n **do**
6: Update f_k by optimizing $L(f_k)$ in Eq. (7)
7: **end for**
8: **end for**
9: **end for**

Complexity Analysis The time complexity of Cool is not a big issue since we can distribute the computation into individual classifiers parallelly. Assume M is the mini-batch size and Λ is the parameter size, then in each mini-batch update, the time complexity for each classifier is $\mathcal{O}(M\Lambda)$. However, for the space complexity, it might be a bottleneck as the storage cost is linearly related to the number of the classifier, i.e., $\mathcal{O}(nM\Lambda)$. Thus, when implementing multiple-network Cool in practices, we have to consider the resource limit.

3. EXPERIMENTS

3.1. Datasets and Baselines

To demonstrate the effectiveness of Cool, we experiment with CIFAR10 and CIFAR100 [16] with pairwise noise [14], asymmetric noise [8], symmetric noise [17] and Clothing1M [18], Food-101N [19], WebVision [20] for real-world noise. We compare Cool with the following two categories of noisy-supervised learning methods. **Single-network Methods:** *Standard*, which directly trains a vanilla classifier on noisy datasets; *Forward* [8], which uses a noise transition matrix for the forward loss correction; *LCCN* [21], which dynamically adjust the transition matrix to safeguard the learning process. We only directly report the result of LCCN on WebVision to save huge labor to reproduce as we adopt the same settings; *Bootstrap* [7], which linearly combines the model predictions and original labels; *MentorNet* [9]. We deploy self-paced MentorNet namely a single model determines small loss samples as useful information and learn with these samples; **Dual-network Methods:** *Decouple* [10], which updates the parameters when the two models disagree; *Co-distillation* [12, 13], which is the dual network version of Bootstrap; *Co-teaching* [14], which is the dual network version of self-paced MentorNet; *Bagging* [22] which takes a vote of multiple classifiers pre-trained on random data partitions. Specifically we use two classifiers.

3.2. Implementation

For CIFAR-10 and CIFAR-100, we follow the same implementation in [14] For real-world datasets, a 50-layer ResNet architecture [3] pre-trained on ImageNet is adopted as the classifier. The images are resized to 256 with respect to shorter sides and then randomly cropped to 224×224 with

random flip, brightness, contrast and saturation. For Clothing1M and Food-101N, the batch size is set to 64 and we run 20 epochs on Clothing1M and 40 epochs on Food-101N. For WebVision, we align the learning settings with Yao et al. [21] to save the labor in reproducing the results.

For all experiments, we use the same architecture for two networks when implementing Decouple, Co-teaching, Co-distillation and Cool as done in [10, 13, 14]. To be fair, all methods use pre-training as warming-up following [8, 13, 14]. Specifically the models are trained as Standard for 4 epochs on Clothing1M and 10 epochs on all other datasets. For dual network methods, two branches are pre-trained separately. For Forward, we use the normalized ground-truth confusion matrix provided in [18] on Clothing1M. For MentorNet and Co-teaching, the noise ratio r is provided as side information as required in [14] for the pre-defined curriculum. For Cool, we set $\lambda = 0.5$ since two networks have the same architecture and are trained in the same manner. We empirically set $\alpha = 0.05$ on CIFAR-10 and Clothing1M while we set $\alpha = 0.1$ on datasets containing much more categories. For the CIFAR datasets, we linearly decrease α since deep models easily fit the small-scale datasets and we fix $\beta = 0.05$.

3.3. Results on CIFAR10 and CIFAR-100

Table 1 summarizes the average test accuracy of dual network Cool and all baselines on CIFAR-10 and CIFAR-100 over the last ten epochs. We can see from the results that dual-network methods such as Co-distillation and Co-teaching generally perform better than their single version namely Bootstrap and MentorNet. Bagging works well on CIFAR-10 but fails on CIFAR-100. The reason is that CIFAR-100 contains more classes thus the data after random partition is insufficient to train a good individual learner. However, adopting similar style of pre-training, Cool manages to achieve the best performance in all of the noise settings. This indicates that Cool can effectively boost two imperfect individual learners. Pointedly for low-level pairwise noise, Cool outperforms the best baseline by 6.24% on CIFAR-10 and 7.24% on CIFAR-100. When r raises to 0.45, all baselines degenerate hard while Cool shows great robustness dealing with high-level noise. According to the results, Cool outperforms the best baseline by 15.88% on CIFAR-10 even impressively surpassing the results of all baselines under low-level noise. For asymmetric noise, Cool achieves the best performances in all settings and manages to outperform the best baseline by 14.47% when $r = 0.45$ on CIFAR-100. Symmetric noise is the hardest noise pattern as the overall test accuracies are low. However, Cool manages to outperform all the baselines.

3.3.1. Counteracting the Memorization Effects

Memorization effects [23] refer to the behavior of DNNs under noisy supervision that the model will firstly learn from

	CIFAR-10						CIFAR-100					
	pairwise		asymmetric		symmetric		pairwise		asymmetric		symmetric	
setting	1	2	3	4	5	6	7	8	9	10	11	12
noise ratio (r)	0.2	0.45	0.2	0.45	0.2	0.5	0.2	0.45	0.2	0.45	0.2	0.5
Standard	76.54	49.73	82.86	69.52	76.71	49.91	50.47	32.51	50.79	31.47	48.63	25.44
Forward	78.06	58.69	83.29	69.59	77.56	51.70	52.67	30.14	54.99	28.41	47.00	25.94
Bootstrap	76.19	49.75	82.90	70.23	76.66	48.83	51.01	31.17	50.35	32.42	47.87	24.74
MentorNet	80.32	58.67	83.62	70.28	80.73	68.70	50.89	32.26	50.27	32.35	52.12	37.96
Decouple	77.30	49.23	83.53	70.24	78.03	48.97	51.28	31.49	50.50	31.72	46.19	23.29
Co-distillation	81.36	51.87	85.28	72.15	81.49	56.76	56.40	34.93	55.27	35.27	54.36	32.41
Co-teaching	83.24	72.74	85.12	76.02	82.09	74.06	54.74	34.08	52.81	34.77	53.84	41.34
Bagging	81.35	57.22	82.58	72.12	79.08	67.90	45.99	26.21	47.24	26.46	43.79	23.98
Cool	89.48	88.62	89.73	82.95	87.88	81.76	63.64	48.38	63.65	49.24	60.36	45.93

Table 1. Average test accuracy (%) on CIFAR-10 & CIFAR-100 over the last ten epochs.

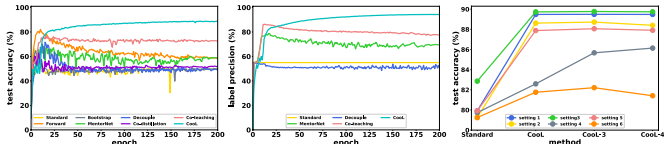


Fig. 1. **Left:** Training curve on CIFAR-10 with pairwise noise ($r = 0.45$); **Middle:** Training curve on CIFAR-10 with pairwise noise ($r = 0.45$); **Right:** Comparison of Cool, Cool-3 and Cool-4.

clean data and eventually fit the noisy labels. This phenomenon can be visualized as rise followed by drop in the test accuracy curve. In the left panel of Figure 1, we trace the test accuracy of all baselines and CoolL under high-level pairwise noise on CIFAR-10. For all baselines, the test accuracy increases at first and then decreases as the training proceeds, which matches the memorization effects. However, the curve of CoolL keeps increasing and then persists at a high level which indicates that the proposed cooperation supervision is reliable enough to counteract the memorization effects.

3.3.2. Reliability of the Supervision

To further assess the reliability of the cooperation supervision in CoolL, We report the label precision which is the ratio of correct supervisions to total supervisions. We compare CoolL with the sample selection methods that leverage different criteria to select the reliable supervision. The label precision curves in setting 2 are depicted in middle panel of Figure 1. We can see that the label precision of CoolL consistently increases with regard to iterations of optimizing and surpasses all the sample selection methods. In the advanced stage of training, the label precision of the cooperation supervision is 94.29% which means CoolL can guarantee the classifiers to learn on a relatively clean dataset. This empirically verifies the reliability of the proposed cooperation supervision.

3.3.3. Two Learners and Beyond

Here we carry out experiments to examine our theoretical findings on the effects of utilizing multiple networks. We

method	Clothing1M	Food-101N
Standard	66.14	75.45
Forward	67.70	75.92
Bootstrap	66.70	74.75
MentorNet	60.00	76.93
Decouple	64.44	72.47
Co-distillation	67.23	77.31
Co-teaching	67.30	77.87
Bagging	67.45	77.14
Cool	70.79	80.94
Cool-3	71.12	81.08

method	acc.@1	acc.@5
Standard	63.11	83.69
Forward	63.10	83.78
LCCN	63.52	84.27
Bootstrap	63.20	83.81
Decouple	61.23	81.53
Co-distillation	63.41	84.14
Co-teaching	-	-
Bagging	63.45	84.19
Cool	63.61	84.32
Cool-3	63.67	84.39

Table 2. Results on Cloth- **Table 3.** Results on WebVi-
ision1M and Food-101N.

implement triple-network CoolL (CoolL-3) and quadruple-network CoolL (CoolL-4) and depict the accuracy along with CoolL and Standard in right panel of Figure 1. We can see from the curves that CoolL-3 generally shows improvement or comparable results with CoolL. This matches our analysis that increasing the number of the classifiers will result in a smaller risk which is closer to the lower bound. CoolL-4 shows improvement in setting 4 while its accuracy slight drops in setting 2 & 6. This is due to insufficient information under the partition of the limited data as we have discussed formerly. As training quadruple classifiers also requires more resources, we may only resort to CoolL or CoolL-3 practically.

3.4. Results on Clothing1M, Food-101N and WebVision

In this section, we empirically verify the effectiveness of CoolL on three large-scale datasets with real-world noise.

For Clothing1M, the results are reported in the left column of Table 2. We can see that dual-network methods generally perform better than the single versions from which they are derived. Although MentorNet does not work well in this setting, Co-teaching manages to surpass Standard by 1.16%. Among all baselines, Forward achieves the best performance with the usage of the ground-truth transition matrix. However CoolL surpasses Forward by 3.09% without using any side information. CoolL-3 further improves CoolL by 0.33% indicating the effectiveness of leveraging multiple classifiers.

For Food-101N, Bootstrap degenerates slightly compared to Standard, while the dual-network version method Co-distillation enjoys a 1.86% gain. Adopting the same small

loss trick, both MentorNet and Co-teaching perform well on Food-101N. MentorNet outperforms Standard by 1.48% and Co-teaching outperforms MentorNet by 0.94% with the use of dual-network structure. Without the knowledge of the ground-truth transition matrix, Forward only improves the test accuracy by 0.47% compared to Standard. Again, CooL manages to outperform the best baseline by 3.07%. Adding another classifier, CooL-3 further obtains a 0.14% gain.

For WebVision, we report both top-1 and top-5 accuracies in Table 3. The results of Co-teaching is vacant due to the absence of the ground-truth noise ratio. We can see from the results that our CooL achieves the best performance and CooL-3 further surpasses CooL. However, the gap between all the methods is trivial which may be on account of the strong open-set noise as suggested in Yao et al. [21].

4. CONCLUSION AND FUTURE WORK

In this paper, we propose a Cooperative Learning paradigm that multiple classifiers work cooperatively with noisy supervision. We demonstrate that our proposed cooperation risk is lower than that associated with individual learners. Then we present a sufficient condition where the risk is negatively correlated to the number of the classifiers. Finally, we introduce the Cooperative Learning framework where the reliable cooperation supervision iteratively boosts the performance of the classifiers. We conduct a range of experiments on the CIFAR datasets to demonstrate the robustness of CooL under synthetic noise and we verify the effectiveness of CooL on three real-world large-scale datasets. We further implement CooL-3 and CooL-4 to show that leveraging more classifiers can have potential gain nonetheless adding more classifiers will consume more resources. Future research directions include finding new means to generate multiple divergent classifiers to achieve lower risk and reducing the parameter space for multiple-network CooL via parameter sharing.

5. REFERENCES

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012.
- [2] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015.
- [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [4] Vikas C Raykar, Shipeng Yu, Linda H Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca Bogoni, and Linda Moy, "Learning from crowds," *Journal of Machine Learning Research*, 2010.
- [5] Rob Fergus, Li Fei-Fei, Pietro Perona, and Andrew Zisserman, "Learning object categories from internet image searches," *Proceedings of the IEEE*, 2010.
- [6] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, "Understanding deep learning requires rethinking generalization," in *ICLR*, 2016.
- [7] Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich, "Training deep neural networks on noisy labels with bootstrapping," in *ICLR*, 2014.
- [8] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu, "Making deep neural networks robust to label noise: A loss correction approach," in *CVPR*, 2017.
- [9] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei, "Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels," in *ICML*, 2018.
- [10] Eran Malach and Shai Shalev-Shwartz, "Decoupling" when to update" from" how to update"," in *NIPS*, 2017.
- [11] Hao Wu, Jiangchao Yao, Jiajie Wang, Yinru Chen, Ya Zhang, and Yanfeng Wang, "Collaborative label correction via entropy thresholding," in *ICDM*, 2019.
- [12] Rohan Anil, Gabriel Pereyra, Alexandre Passos, Robert Ormandi, George E. Dahl, and Geoffrey E. Hinton, "Large scale distributed neural network training through online distillation," in *ICLR*, 2018.
- [13] Guocong Song and Wei Chai, "Collaborative learning for deep neural networks," in *NIPS*, 2018.
- [14] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," in *NIPS*, 2018.
- [15] Yuncheng Li, Jianchao Yang, Yale Song, Liangliang Cao, Jiebo Luo, and Li-Jia Li, "Learning from noisy labels with distillation," in *ICCV*, 2017.
- [16] Alex Krizhevsky and Geoffrey Hinton, "Learning multiple layers of features from tiny images," Tech. Rep., Citeseer, 2009.
- [17] Brendan Van Rooyen, Aditya Menon, and Robert C Williamson, "Learning with symmetric label noise: The importance of being unhinged," in *NIPS*, 2015.
- [18] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang, "Learning from massive noisy labeled data for image classification," in *CVPR*, 2015.
- [19] Kuang-Huei Lee, Xiaodong He, Lei Zhang, and Linjun Yang, "Cleannet: Transfer learning for scalable image classifier training with label noise," in *CVPR*, 2018.
- [20] Wen Li, Limin Wang, Wei Li, Eirikur Agustsson, and Luc Van Gool, "Webvision database: Visual learning and understanding from web data," *arXiv:1708.02862*, 2017.
- [21] Jiangchao Yao, Hao Wu, Ya Zhang, Ivor W Tsang, and Jun Sun, "Safeguarded dynamic label regression for noisy supervision," *AAAI*, 2019.
- [22] Leo Breiman, "Bagging predictors," *Machine learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [23] Devansh Arpit, Stanisław Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al., "A closer look at memorization in deep networks," in *ICML*, 2017.