

Deep Discriminative CNN with Temporal Ensembling for Ambiguously-Labeled Image Classification

Yao Yao¹, Chen Gong¹, Jiehui Deng¹, Xiuhua Chen¹, Jianxin Wu³ and Jian Yang^{1,2*}

AAAI 2020





□ The groundtruth of each training image is unique and should be available in training phase.

Unfortunately, the images may lack clear labels and manually labeling them will incur unaffordable monetary or time cost.



Partial Label Learning

Newsletter



Barcelona news: Lionel Messi and Luis Suarez hold private talks over five player concerns

Crowdsourcing



Annotator 1: Horse Annotator 2: Donkey Annotator 3: Mule

Generally, it is a multiple classification problem.

□ The images are associated with multiple candidate labels, and only one of them is valid

Formal Definition



\square Input space: $\mathcal{X} \in \mathbb{R}^d$ is *d*-dimensional.

 \square Output space: $\mathcal{Y} = \{1, 2, \dots, c\}$ includes c classes.

 \square Training set: $\mathcal{D} = \{(\mathbf{x}_i, S_i) | 1 \leq i \leq n\}$ has n ambiguously labeled examples.

False positive labels Candidate label set: $S_i = A_i \cup \{y_i\}$

Groundtruth label

 \square Target: train a classifier $f: \mathcal{X} \to \mathcal{Y}$ from the training set \mathcal{D} so that the correct predictions can be made on test examples.

Motivations



- To solve this problem, the common strategy is to disambiguate the set of candidates of each example.
- However, existing methods are usually short of representation ability and discrimination ability.

	Representation ability	Discrimination ability
Reasons	Shallow learning frameworks	Imperfect disambiguation techniques.
Contributions	Employ deep convolutional neural network	Employ entropy minimization regularizer and temporal ensembling technique

Methods





Loss function:

Discrimination term

$$Loss = \mathcal{L}_f(\mathbf{Y}, \hat{\mathbf{Y}}) + \alpha \mathcal{L}_d(\hat{\mathbf{Y}}) + T(t) \cdot \mathcal{L}_t(\bar{\mathbf{Y}}, \hat{\mathbf{Y}})$$

Fidelity term

Temporal Ensembling Term



The predicted probability of the labels in non-candidate label set should be zero.

$$\mathcal{L}_f(\mathbf{Y}, \hat{\mathbf{Y}}) = -\frac{1}{n} \sum_{i=1}^n (\mathbf{1} - \mathbf{y}_i)^\top \log(\mathbf{1} - \hat{\mathbf{y}}_i)$$

- \blacktriangleright Label matrix: $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \{0, 1\}^{n imes c}$
 - $y_{ij} = 1$: if the *j*-th label is a candidate label of the image \mathbf{X}_i
 - $y_{ij} = 0$: otherwise
- \blacktriangleright Predicted matrix: $\hat{\mathbf{Y}} = [\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_n] \in [0, 1]^{n \times c}$
 - \hat{y}_{ij} : the probability of the image \mathbf{X}_i belongs to the *j*-th class

•
$$\hat{\mathbf{y}}_{ij} \ge 0$$
 and $\sum_{j=1}^{c} \hat{\mathbf{y}}_{ij} = 1$

Discrimination Term



Minimizing entropy makes the potential groundtruth label become prominent among all labels.





Temporal Ensembling Term



Assemble the model predictions of different epochs and regard them as the auxiliary supervision information for the next epoch.

$$\mathcal{L}_t(\bar{\mathbf{Y}}, \hat{\mathbf{Y}}) = -\frac{1}{n} \sum_{i=1}^n w_i \cdot \bar{\mathbf{y}}_i^\top \log \hat{\mathbf{y}}_i$$

$$\mathbf{S}^{(t)} = \gamma \mathbf{S}^{(t-1)} + (1-\gamma) \hat{\mathbf{Y}}^{(t)}$$

Assembled predictions at t-th epoch

$$\bar{\mathbf{Y}}^{(t+1)} = \mathbf{S}^{(t)} / (1 - \gamma^t)$$

Training target of (t+1)-th epoch

$$w_i = \begin{cases} 0 & m_i \le 0\\ {m_i}^2 & otherwise \end{cases}$$

$$m_i = \max_{y_j \in S_i} \overline{y}_{ij} - \max_{y_k \notin S_i} \overline{y}_{ik}$$

Formation of T(t)



 \succ T(t) is a time-dependent weighting function.

> Roughly speaking, T(t) increases with the number of epochs.

$$T(t) = T_{max} \exp\left[-5(1-t)^2\right]$$

> T_{max} denotes the maximum value that T(t) could reach.

- \succ t increases linearly from zero to one during the rising phase.
- At the initial training phase, the network mainly learns from the original ambiguous labels.
- Then it gradually learns from the assembled predictions when then training process proceeds

Datasets



- > Synthesized datasets:
 - Fashion MNIST

- Real-world datasets:
 - Yahoo!News

• SVHN

• Lost

Datasets	# Images	# Classes	Avg # labels
FM-v1	70,000	10	2
FM-v3	70,000	10	4
SVHN-v1	99,289	10	2
SVHN-v3	99,289	10	4
Lost	1,122	16	2.23
Yahoo!News	14,322	38	1.44

Characteristics of the adopted datasets.



Experiments on Synthesized Datasets

	FM-v1	FM-v3	SVHN-v1	SVHN-v3
RegISL	-	-	-	-
SURE	-	-	-	-
WMCar-ICE	-	-	-	-
MCar	0.913 ± 0.003 $ullet$	0.825 ± 0.002 \bullet	$0.796 \pm 0.004 \bullet$	$0.545\pm0.004\bullet$
PLKNN	0.897 ± 0.004 $ullet$	0.848 ± 0.004 $ullet$	$0.736\pm0.003 \bullet$	$0.636\pm0.002\bullet$
M3PL	0.884 ± 0.002 $ullet$	0.874 ± 0.004 $ullet$	$0.827\pm0.002\bullet$	0.788 ± 0.003 $ullet$
IPAL	0.912 ± 0.005 \bullet	0.905 ± 0.003 $ullet$	$0.798 \pm 0.003 \bullet$	$0.777\pm0.001\bullet$
DCNN	0.902 ± 0.007 $ullet$	0.890 ± 0.008 $ullet$	$0.922\pm0.003\bullet$	$0.908\pm0.007\bullet$
D^2CNN	$\textbf{0.936} \pm \textbf{0.002}$	$\textbf{0.927} \pm \textbf{0.003}$	$\textbf{0.937} \pm \textbf{0.003}$	$\textbf{0.929} \pm \textbf{0.001}$

- D²CNN achieves superior performance against other baselines on these synthesized datasets.
- The performances of all baselines decrease, when the number of the candidate labels increase.



	Lost	Yahoo!News
RegISL	0.761 ± 0.037 $ullet$	$0.598\pm0.016ullet$
SURE	0.794 ± 0.037 $ullet$	$0.729\pm0.010ullet$
WMCar-ICE	0.795 ± 0.020 $ullet$	$0.705\pm0.010ullet$
MCar	0.743 ± 0.011 $ullet$	0.671 ± 0.010 \bullet
PLKNN	0.651 ± 0.012 $ullet$	$0.562\pm0.017ullet$
M3PL	0.678 ± 0.032 $ullet$	0.613 ± 0.001 $ullet$
IPAL	0.790 ± 0.034 $ullet$	$0.647\pm0.017ullet$
DCNN	0.580 ± 0.031 $ullet$	0.740 ± 0.006 $ullet$
D^2CNN	$\textbf{0.838} \pm \textbf{0.014}$	$\textbf{0.833} \pm \textbf{0.009}$

- The accuracy of D²CNN is higher than the second best method on *Lost* dataset by approximately 4%.
- On *Yahho!News* dataset, D²CNN significantly outperforms other baselines and leads the second best method with the margin of 9.3%.

The robustness to proportion of ambiguously-labeled images





Fashion-Mnist dataset when r = 3





Algorithm Analysis



The effectiveness of \mathcal{L}_d and \mathcal{L}_t

The effectiveness of varied T(t)

THANKS