



- •The key of generic object, scene or action recognition is the separable of features. Because the classes are constant.
- •The deeply learned features are required to be generalized enough for identifying new unseen classes without label prediction.

3. Formula of contrastive-center loss



- rate.
- inevitable in contrastive loss and triple loss.
- but also image classification.

7. Experiments on MNIST and Cifar-10									
• MNIST:		stage 1		stage 2		stage	3	stage 4	
(use a small network	Layer	conv	pool	conv	pool	conv	pool	FC	
for visualization and	LeNets	$(5,20)_{/1,0}$	$2_{/2,0}$	$(5, 50)_{/1,0}$	$2_{/2,0}$			500	
accuracy comparison)	LeNets++	$(5, 32)_{/1,2} \times 2$	$2_{/2,0}$	$(5, 64)_{/1,2} \times 2$	$2_{/2,0}$	$(5, 128)_{/1,2} \times$	$< 2 2_{/2,0}$	2	
	Method Softmax Center loss		Accuracy(% 98.8 98.94 99.17						
• Cifar-10: (use 20-layer ResNet)	Cifar-10: e 20-layer ResNet)MethodAccuracy(%)20-layer ResNet [5]91.2520-layer ResNet(our implementation based on center loss [13])92.1							7(%)	
	20-layer ResNet(our contrastive-center loss)						92.45		

CONTRASTIVE-CENTER LOSS FOR DEEP NEURAL NETWORKS

Ce Qi, Fei Su School of Information and Communication Engineering Beijing University of Posts and Telecommunications, Beijing

1. Introduction

•For face recognition task, the deeply learned features need to be not only separable but also discriminative. Because the classes of faces are inconstant.

L, L_s , L_c , L_{ct-c} denote the total loss, softmax loss, center loss and contrastive-center

 λ denotes the scalar used for balancing the two loss functions.

m denotes the number of training samples in a mini-batch.

 $x_i \in \mathbf{R}^d$ denotes the *i*th training sample with dimension d.

d is the feature dimension.

 $\boldsymbol{y_i}$ denotes the label of x_i .

 $c_{y_i} \in \mathbf{R}^d$ denotes the y_i th class center of deep features with dimension d.

k denotes the number of class.

 $oldsymbol{\delta}$ is a constant used for preventing the denominator equal to 0.

In our experiments, we set $\delta = 1$ by default.

5. Details of contrastive-center loss

•Centers are updated based on mini-batch with an adjustable learning

•Contrastive-center loss enjoys the same requirement as the softmax loss and needs no complex sample mining and recombination, which is

•Contrastive-center loss does help on tasks of not only face recognition

•In general, for tasks using class labels, contrastive-center can do help.

2. Discriminative Feature Learning

- •SOFTMAX LOSS: encourages the separability of features. •CENTER LOSS: learns a center for deep features of each class, but only penalizes the distances between the deep features and their corresponding class centers.
- and inter-class separability, by penalizing the contrastive values between: (1)the distances of training samples to their corresponding class centers, and (2)the sum of the distances of training samples to their non-corresponding class centers.

4. Details of contrastive-center loss About code optimazation on GPU: • Derivative: $\left\|x_{i}-c_{y_{i}}\right\|^{2} \sum_{j=1, j\neq y_{i}}^{k} (x_{i}-c_{j})$ $\frac{\partial L_{ct-c}}{\partial x_i}$ $\overline{\left(\sum_{j=1, j\neq y_i}^k \|x_i - c_j\|^2\right)} + \delta$ $\left[\left(\sum_{j=1,j\neq y_{i}}^{k}\left\|x_{i}-c_{j}\right\|^{2}\right)+\delta\right]^{2}$ 128*512 kernels , while the latter uses 512*100000 kernels. $\frac{1}{\|y_i\|^2} = n$ $\left(\sum_{j=1,j\neq y_i}^k \|x_i - c_j\|\right)$ • $\frac{\partial L_{ct-c}}{\partial c_n} = \sum_{i=1}^m$ $(x_i - c_n) \left\| x_i - c_{y_i} \right\|$ $\left\| \left[\sum_{j=1, j \neq y_i}^{k} \|x_i - c_j\|^2 + \delta \right]^2 if y_i \neq n \right]$ time

6. Visualization on MNIST

- Softmax loss: features are separable, but not discriminative.
- Softmax loss+center loss: features are separable and have better intra-class compactness.
- Softmax loss + contrastive-center loss: features are separable and have better intraclass compactness and inter-class separability $(30 \text{ vs } 10^{-15}, 2^{-3} \text{ times than center loss}).$

8. Experiments on LFW



•CONTRASTIVE-CENTER LOSS: simultaneously considers intra-class compactness

- For 128,512,10 vs 128,512,100000 (batch size, feature dimension and number of class), the former uses

Keep the variables with high time complexity used in forward and backward in memory to save computing



(c) our contrastive-center loss

	Images	Networks	$A_{\rm ccuracy}(\%)$
	images		
	4M	3	97.35
	10M	5	98.37
	0.5M	1	98.60
	0.5M	1	98.62
		1	98.70
	2.6M	1	98.95
	0.494, 414M	1	97.73
	0.494, 414M	1	98.43
	0.455, 594M	1	97.47
)	0.494, 414M	1	98.43
	0.455, 594M	1	98.55
	0.455,594M	1	98.68