

Large-Margin Softmax loss for Convolutional Neural Networks

M2019170
허성실



0

Abstract

- Cross-entropy loss together with softmax is one of the most common used super-vision components
- But the component does not explicitly encourage discriminative learning of features
- L-Softmax loss explicitly encourages **Intra-class compactness and inter-class separability**
- Not only can adjust the **desired margin** but also can **avoid overfitting**
- Can be optimized by typical stochastic gradient descent



1

Introduction

- CNN needs more **discriminative** information
- **Intra-class compactness and inter-class separability** should simultaneously maximized.
- Contrastive loss and triplet loss solved that problem, but needs too many training samples
- They define the softmax loss as the combination of a cross-entropy loss, a softmax function and the las fully connected layer
- The purpose of this paper is **to generalize** the softmax loss to a more general large-marin softmax loss
- This is done by incorporating a preset constant m multiply with the angle between sample

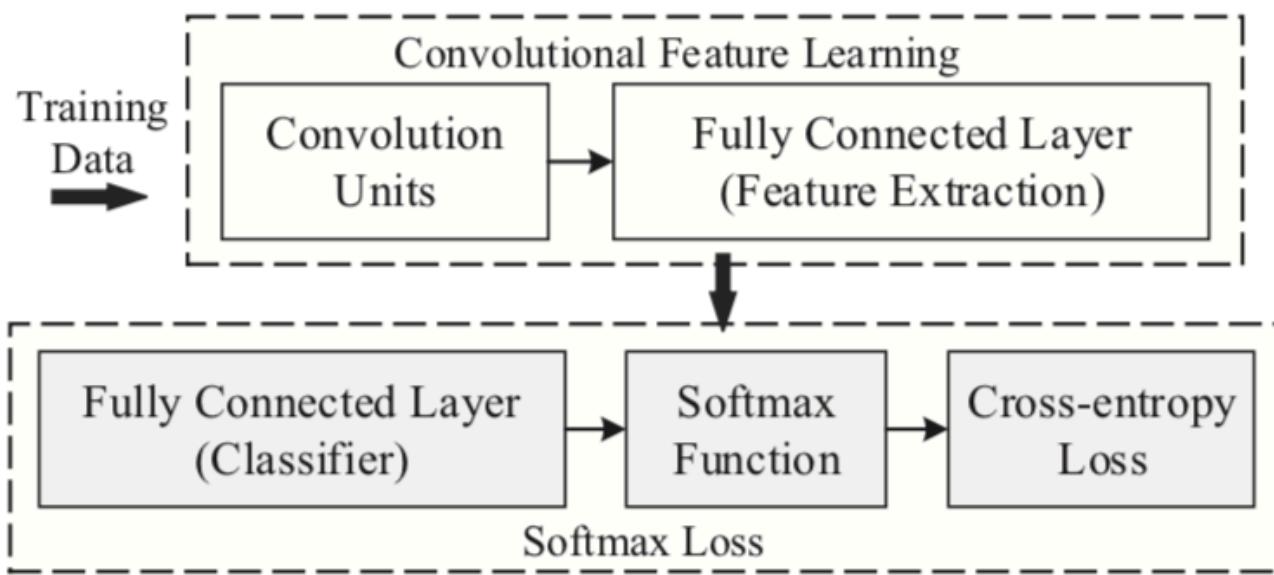


Figure 1. Standard CNNs can be viewed as convolutional feature learning machines that are supervised by the softmax loss.



2

Related Work

- **Softmax loss + contrastive loss**
 - > need pair of training samples
 - > same class requires having similar features
- **Triplet loss**
 - > needs 3 training samples
 - > minimize. The distance between the anchor and positive samples
- They both need a **carefully designed pair selection procedure**

- Softmax loss

$$L = \frac{1}{N} \sum_i L_i = \frac{1}{N} \sum_i -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$

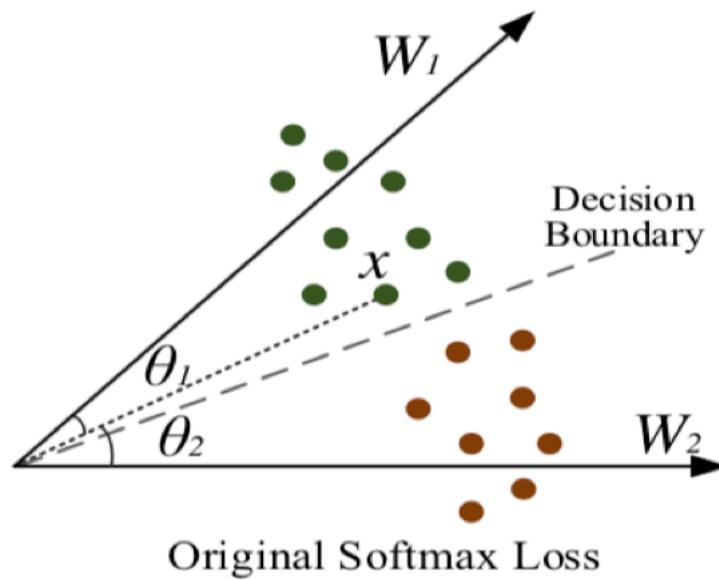
||

$$L_i = -\log \left(\frac{e^{\|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\| \cos(\theta_{y_i})}}{\sum_j e^{\|\mathbf{W}_j\| \|\mathbf{x}_i\| \cos(\theta_j)}} \right)$$

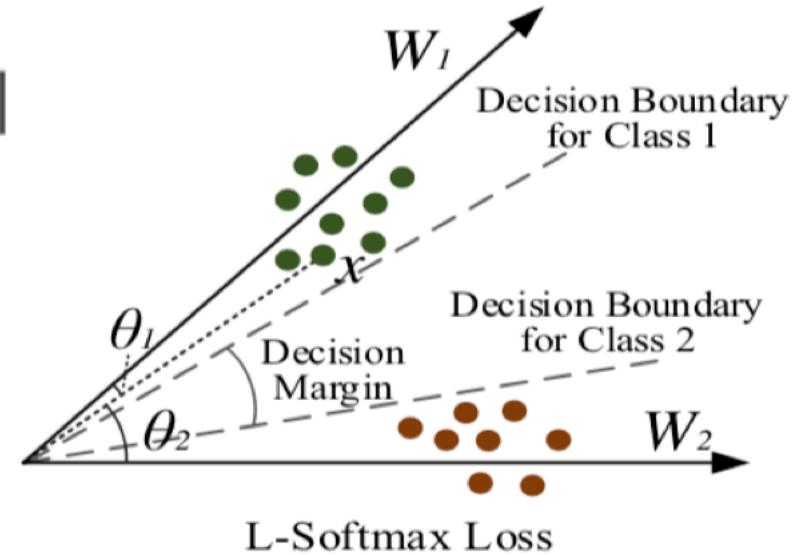


3

Large-Margin Softmax Loss



$$\|\mathbf{W}_1\| = \|\mathbf{W}_2\|$$



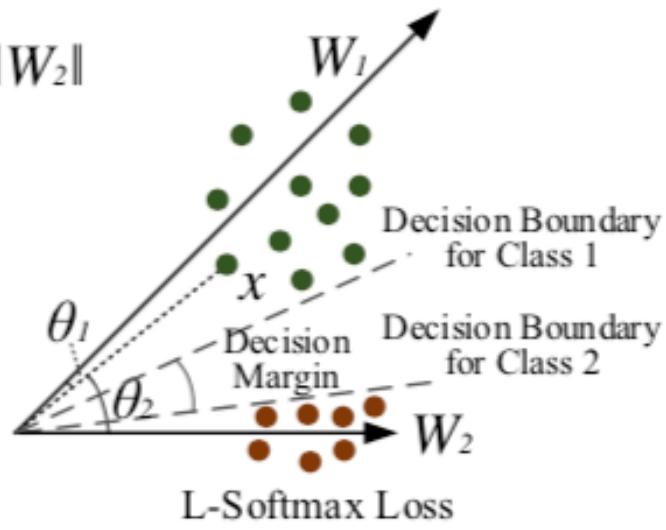
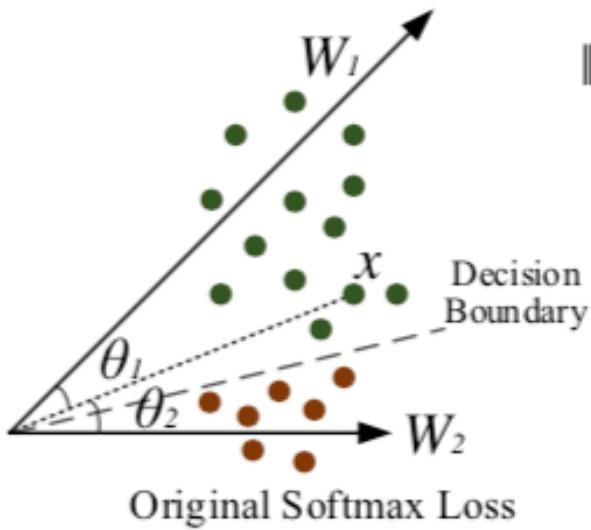
$$\mathbf{W}_1^T \mathbf{x} > \mathbf{W}_2^T \mathbf{x}$$

$$\|\mathbf{W}_1\| \|\mathbf{x}\| \cos(\theta_1) > \|\mathbf{W}_2\| \|\mathbf{x}\| \cos(\theta_2)$$

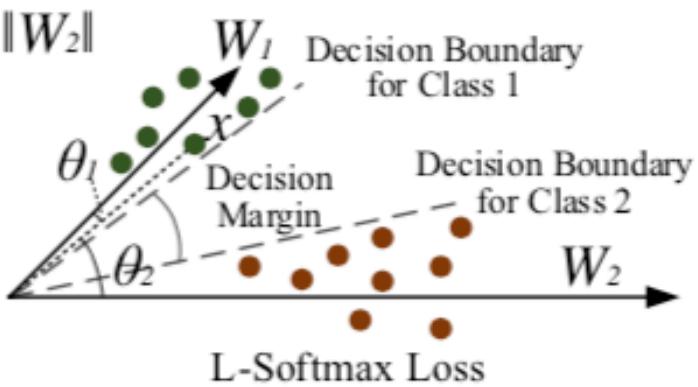
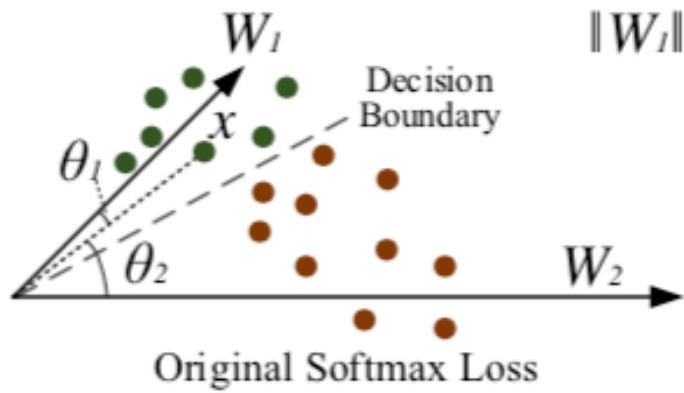
$$\|\mathbf{W}_1\| \|\mathbf{x}\| \cos(m\theta_1) > \|\mathbf{W}_2\| \|\mathbf{x}\| \cos(\theta_2) \quad (0 \leq \theta_1 \leq \frac{\pi}{m})$$

$$\begin{aligned} \|\mathbf{W}_1\| \|\mathbf{x}\| \cos(\theta_1) &\geq \|\mathbf{W}_1\| \|\mathbf{x}\| \cos(m\theta_1) \\ &> \|\mathbf{W}_2\| \|\mathbf{x}\| \cos(\theta_2). \end{aligned}$$

- The original softmax loss requires $\theta_1 < \theta_2$ to classify the sample \mathbf{x} as class 1, while the L-Softmax loss requires $m\theta_1 < \theta_2$ to make the same decision.



* The larger W is, the larger the feasible angle of its corresponding class is.



L-softmax

$$L_{softmax} = -\log \frac{\exp(W_y^T \mathbf{x})}{\sum_{j=1}^C \exp(W_j^T \mathbf{x})} = -\log \frac{\exp(\|W_y\| \|\mathbf{x}\| \cos(\theta_y))}{\sum_{j=1}^C \exp(\|W_j\| \|\mathbf{x}\| \cos(\theta_j))}$$

$$L_{L-softmax} = -\log \frac{\exp(\|W_y\| \|\mathbf{x}\| \psi(\theta_y))}{\exp(\|W_y\| \|\mathbf{x}\| \psi(\theta_y)) + \sum_{j \neq y} \exp(\|W_j\| \|\mathbf{x}\| \cos(\theta_j))}$$

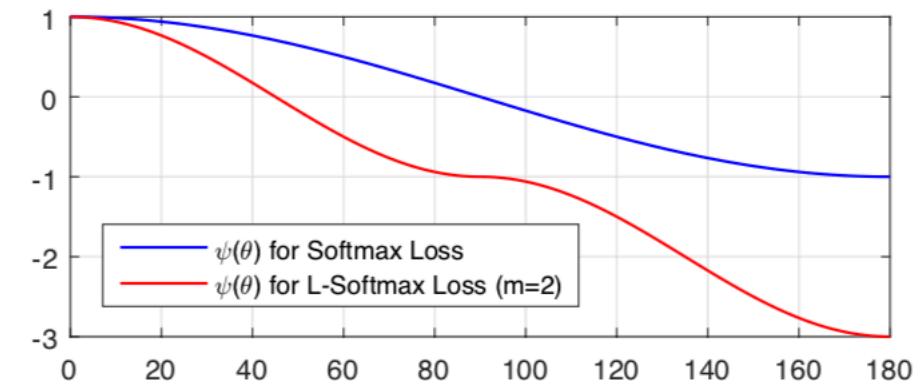


Figure 3. $\psi(\theta)$ for softmax loss and L-Softmax loss.

$$\psi(\theta) = \begin{cases} \cos(m\theta), & 0 \leq \theta \leq \frac{\pi}{m} \\ \mathcal{D}(\theta), & \frac{\pi}{m} < \theta \leq \pi \end{cases}$$

$$\psi(\theta) = (-1)^k \cos(m\theta) - 2k, \quad \theta \in [\frac{k\pi}{m}, \frac{(k+1)\pi}{m}]$$

- Transform $\cos(m\theta)$ into combinations of $\cos(\theta)$:

$$\begin{aligned}\cos(m\theta_{y_i}) &= C_m^0 \cos^m(\theta_{y_i}) - C_m^2 \cos^{m-2}(\theta_{y_i})(1 - \cos^2(\theta_{y_i})) \\ &\quad + C_m^4 \cos^{m-4}(\theta_{y_i})(1 - \cos^2(\theta_{y_i}))^2 + \dots \\ &\quad (-1)^n C_m^{2n} \cos^{m-2n}(\theta_{y_i})(1 - \cos^2(\theta_{y_i}))^n + \dots\end{aligned}$$

- Represent $\cos(\theta)$ as

$$\frac{\mathbf{W}_j^T \mathbf{x}_i}{\|\mathbf{W}_j\| \|\mathbf{x}_i\|}$$

$$\begin{aligned}\cos(2\theta) &= 2\cos^2\theta - 1. \\ C_2 \cos^2(\theta) &\sim C_2^2 \cos^4\theta (1 - \cos^2\theta) \\ &= C_2^2 \cos^2\theta - C_2^2 (1 - \cos^2\theta) \\ &= \cancel{C_2^2} \cos^2\theta - 1 + \cos^2\theta \\ &= 2\cos^2\theta - 1\end{aligned}$$

$$\begin{aligned}\cos(4\theta) &= 2\cos^2(2\theta) - 1 \\ &= 2(\cos(2\theta) \cdot \cos(2\theta)) - 1 \\ &= 2(2\cos^2\theta - 1)^2 - 1 \\ &= 2(4\cos^4\theta - 4\cos^2\theta + 1) - 1 \\ &= 8\cos^4\theta - 8\cos^2\theta + 1 \\ C_4^0 \cos^4(\theta) &- C_4^2 \cos^2\theta (1 - \cos^2\theta) + C_4^4 \cos^2\theta (1 - \cos^2\theta)^2 \\ &= \cos^4\theta - 6(\cos^2\theta - \cos^4\theta) + (1 - 2\cos^2\theta + \cos^4\theta) \\ &= \cancel{\cos^4\theta} - 6\cos^2\theta + \cancel{6\cos^4\theta} + 1 - 2\cos^2\theta + \cancel{\cos^4\theta} \\ &= 8\cos^4\theta - 8\cos^2\theta + 1\end{aligned}$$

- The only difference between the original loss and the L-softmax loss lies in f_{y_i}

$$\begin{aligned}
 f_{y_i} &= (-1)^k \cdot \|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\| \cos(m\theta_i) - 2k \cdot \|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\| \\
 &= (-1)^k \cdot \|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\| \left(C_m^0 \left(\frac{\mathbf{W}_{y_i}^T \mathbf{x}_i}{\|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\|} \right)^m - \right. \\
 &\quad \left. C_m^2 \left(\frac{\mathbf{W}_{y_i}^T \mathbf{x}_i}{\|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\|} \right)^{m-2} (1 - \left(\frac{\mathbf{W}_{y_i}^T \mathbf{x}_i}{\|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\|} \right)^2) + \dots \right) \\
 &\quad - 2k \cdot \|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\|
 \end{aligned}$$

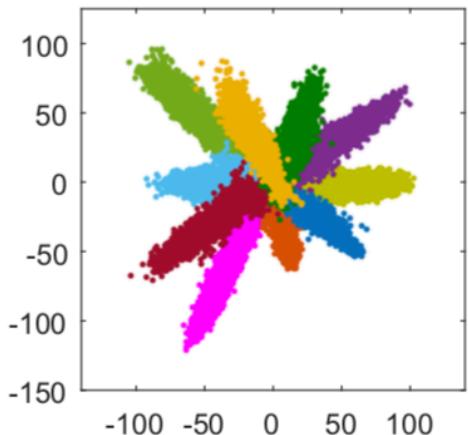
..

- In practice, they seek to minimize:

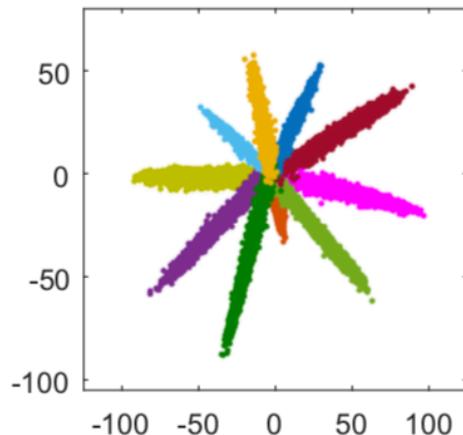
$$f_{y_i} = \frac{\lambda \|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\| \cos(\theta_{y_i}) + \|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\| \psi(\theta_{y_i})}{1 + \lambda}$$

- Start with large $\lambda(0.1)$ and gradually reduce to a very small value.(0.01, 0.001...)

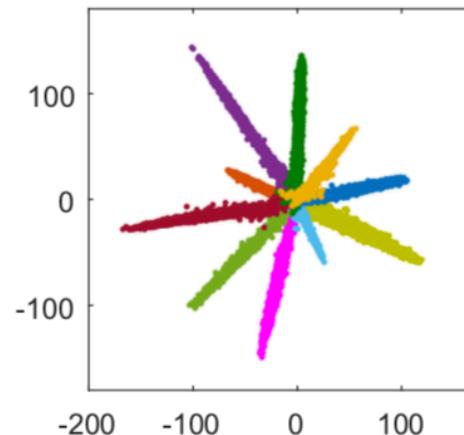
MNIST example



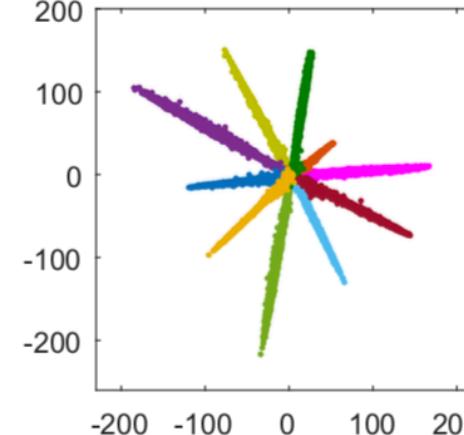
Training Set ($m=1$, Softmax)



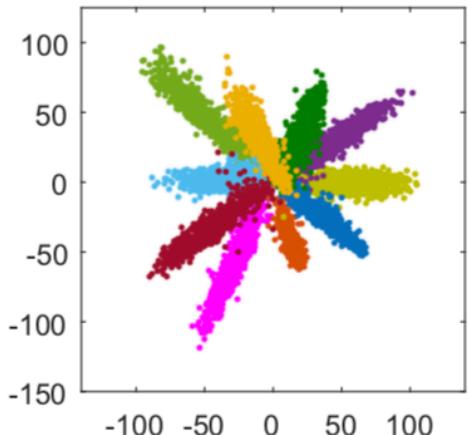
Training Set ($m=2$)



Training Set ($m=3$)

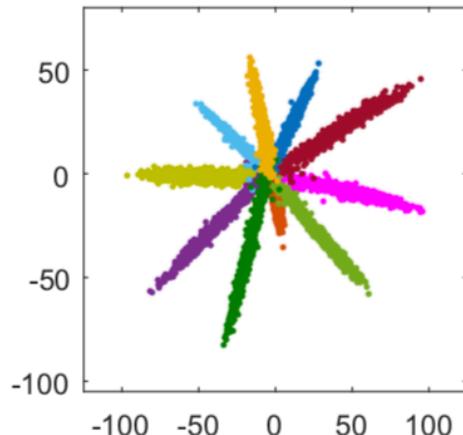


Training Set ($m=4$)



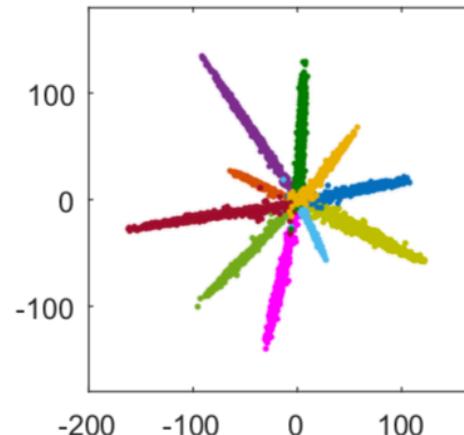
Testing Set ($m=1$, Softmax)

Testing Accuracy: 98.45%



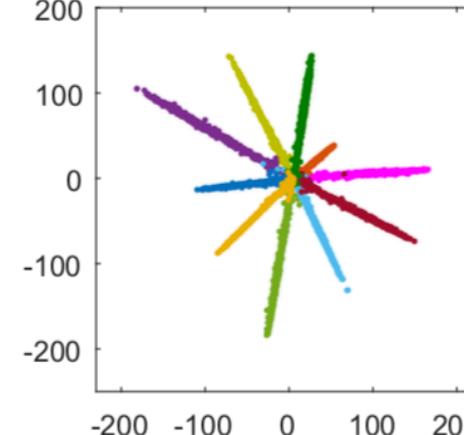
Testing Set ($m=2$)

Testing Accuracy: 98.96%



Testing Set ($m=3$)

Testing Accuracy: 99.22%



Testing Set ($m=4$)

Testing Accuracy: 99.34%



Less overfitting problem compared to original Softmax

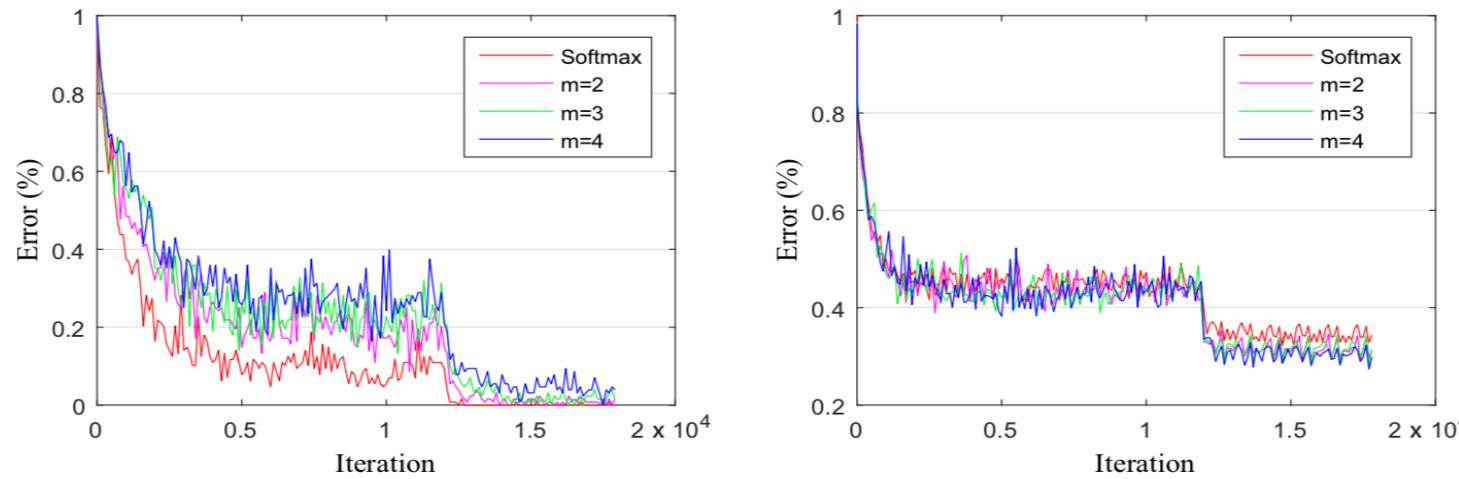


Figure 6. Error vs. iteration with different value of m on CIFAR100. The left shows training error and the right shows testing error.

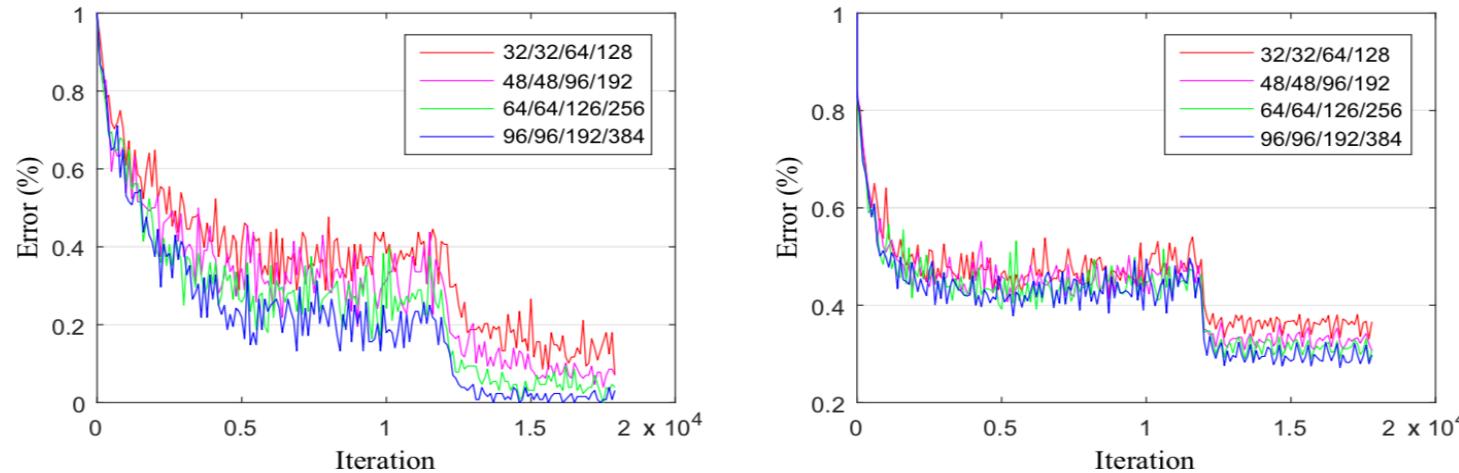


Figure 7. Error vs. iteration ($m=4$) with different number of filters on CIFAR100. The left (right) presents training (testing) error.

Method	Outside Data	Accuracy
FaceNet (Schroff et al., 2015)	200M*	99.65
Deep FR (Parkhi et al., 2015)	2.6M	98.95
DeepID2+ (Sun et al., 2015)	300K*	98.70
(Yi et al., 2014)	WebFace	97.73
(Ding & Tao, 2015)	WebFace	98.43
Softmax	WebFace	96.53
Softmax + Contrastive	WebFace	97.31
L-Softmax (m=2)	WebFace	97.81
L-Softmax (m=3)	WebFace	98.27
L-Softmax (m=4)	WebFace	98.71

Table 5. Verification performance (%) on LFW dataset. * denotes the outside data is private (not publicly available).



5

Conclusion

- Adjustable margin with parameter m
- Clear intuition and geometric interpretation
- Show clear advantages over current state-of-the-art CNNs