

CosFace: Large Margin Cosine Loss for Deep Face Recognition

Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Zhifeng Li, Dihong Gong, Jingchao Zhou, Wei Liu
Tencent AI Lab, China

{hawelwang, yitongwang, encorezhou, denisji, michaelzfli}@tencent.com
gongdihong@gmail.com sagazhou@tencent.com wliu@ee.columbia.edu

Abstract

Face recognition has achieved revolutionary advancement owing to the advancement of the deep convolutional neural network (CNN). The central task of face recognition, including face verification and identification, involves face feature discrimination. However, traditional softmax loss of deep CNN usually lacks the power of discrimination. To address this problem, recently several loss functions such as central loss [1], large margin softmax loss [2], and angular softmax loss [3] have been proposed. All these improvement algorithms share the same idea: maximizing inter-class variance and minimizing intra-class variance. In this paper, we design a novel loss function, namely large margin cosine loss (LMCL), to realize this idea from a different perspective. More specifically, we reformulate the softmax loss as cosine loss by L2 normalizing both features and weight vectors to remove radial variation, based on which a cosine margin term m is introduced to further maximize decision margin in angular space. As a result, minimum intra-class variance and maximum inter-class variance are achieved by normalization and cosine decision margin maximization. We refer to our model trained with LMCL as CosFace. To test our approach, extensive experimental evaluations are conducted on the most popular public-domain face recognition datasets such as MegaFace Challenge, Youtube Faces (YTF) and Labeled Face in the Wild (LFW). We achieve the state-of-the-art performance on these benchmark experiments, which confirms the effectiveness of our approach.

1

1. Introduction

Recently progress on the development of deep convolutional neural network (CNN) [4, 5, 6, 7, 8] has significantly improved the state-of-the-art performance for a wide variety of computer vision tasks, which makes deep CNN a domi-

¹This paper is under peer-review as a conference paper in 2017. The results have been published on Megaface official website (<http://megaface.cs.washington.edu/>) in Sept. 2017

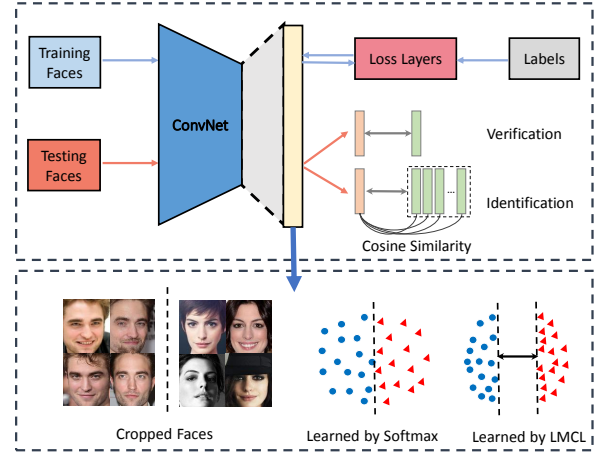


Figure 1. An overview of the proposed CosFace framework. At training phase, the discriminative face features are learned with large margin between different classes. At testing phase, the testing data feed into CosFace to extract face features which are later used to compute the cosine similarity score to perform face verification and identification.

nant machine learning approach for computer vision. Face recognition, as one of the most common computer vision tasks, usually include two sub-tasks: face verification and face identification, where verification performs the comparison on a pair of faces to decide whether they derive from the same subject, while face identification identifies a face from a gallery of faces. Both of these two tasks involve three stages: face detection, feature extraction, and classification. Deep CNN is able to extract clean high-level features, which makes it possible to achieve superior performance with relatively simple classification network: usually, multilayer perceptron networks followed by softmax [9, 10]. However, recent research studies [1, 2, 3] found that traditional softmax is insufficient to maximize the discrimination power for classification.

To encourage better discrimination performance, many research studies have been carried out [1, 11, 12, 13, 14, 3]. All these studies share the same idea for maximum discrimi-

nation capability: maximizing inter-class variance and minimizing intra-class variance. For example, [1, 11, 12, 13, 14] propose to adopt multi-loss learning in order to increase the feature discrimination power. While these methods improve classification performance over traditional softmax loss, they usually come with some extra limitations. For [1], it only explicitly minimizes intra-class variance while ignoring the inter-class variances, which may result in suboptimal solutions. For [11, 12, 13, 14], they require thoroughly scheming the mining of pair or triplet samples, which is an extremely time-consuming procedure. Very recently, [3] have been proposed to address this problem from a different perspective. More specifically, [3] (A-softmax) introduces angular margin for larger inter-class variance. A-softmax projects the original Euclidean space of features to an angular space by constraining the norm of weight to 1 with zero bias.

Compared to Euclidean margin suggested by [1, 11, 13], the angular margin is preferred because the cosine of the angle has intrinsic consistency with softmax. However, on a further note, it seems more natural to directly introduce the cosine margin between different classes. Moreover, the formulation of cosine matches the similarity measurement that is frequently applied to face recognition. From above perspectives, cosine margin drives a straightforward way to improve cosine-related discriminative information better than Euclidean margin or angular margin.

In this paper, we reformulate the softmax loss as cosine loss by L2 normalizing both features and weight vectors to remove radial variation, based on which a cosine margin term m is introduced to further maximize decision margin in cosine angular space. Specifically, we invent an ingenious algorithm named Large Margin Cosine Loss (LMCL) which takes the normalized features as input to learn highly discriminative features by maximizing inter-class cosine margin. Formally, we define a margin m such that the decision boundary is given by $\cos(\theta_1) - m = \cos(\theta_2)$, where θ_i is the angle between the feature and weight of class i .

For comparison, the decision boundary of the A-softmax is defined over angular space by $\cos(m\theta_1) = \cos(\theta_2)$, which has difficulty in optimization due to the non-monotonicity of the cosine function. To overcome such difficulty, one has to employ extra facility with a ad-hoc piecewise function for the A-softmax. More importantly, the decision margin of the A-softmax depends on θ , which leads to different margins for different classes. As a result, in the decision space, some inter-class features have larger margin while others have the smaller margin, which causes reduction of discrimination power. Unlike the A-softmax, our approach defines decision margin in the cosine space, which avoids the aforementioned shortcomings.

Based on the LMCL, we develop a sophisticated deep model called CosFace, as shown in figure 1. At training

phase, LMCL guides the ConvNet to learn features with large cosine margin. At testing phase, the face features are extracted from ConvNet to perform either face verification or face identification. We summarize the contribution as follows:

- (1) We embrace the idea of maximizing inter-class variation and minimizing intra-class variation and propose a novel loss, called LMCL, to learn highly discriminative deep features for face recognition.
- (2) We provide a reasonable theoretical analysis based on the hyperspherical feature distribution encouraged by LMCL.
- (3) The proposed approach improves the state-of-the-art performance over most of the benchmarks on popular face databases including LFW[15], YTF[16] and Megaface[17, 18].

2. Related Work

Deep Face Recognition. Recently face recognition has achieved significant progress thanks to the great success of deep CNN models [5, 4, 19, 7]. In DeepFace [9] and DeepID [10], face recognition is treated as a multi-class classification problem and deep CNN models are first introduced to learn features on large multi-identities datasets. DeepID2 [20] employs identification signals and verification signals to achieve better feature embedding. Recent works DeepID2+ [21] and DeepID3 [22] further explore the advanced network structures to boost performance. FaceNet [23] uses triplet loss to learn a Euclidean space embedding and a deep CNN is trained on nearly 200 million face images, which leads to the state-of-the-art performance. Other approaches [24, 25] also prove effectiveness of the deep CNN on face recognition.

Loss Functions. Loss function plays an important role on deep feature learning. Contrastive loss [11, 12] and triplet loss [13, 14] are usually used to increase Euclidean margin for better feature embedding. Wen et al. [1] propose center loss to learn centers for deep features of each identity and use the centers to reduce intra-class variance. Liu et al. [2] propose large margin softmax (L-Softmax) by adding angular constraints to each identity to improve feature discrimination. Angular softmax (A-Softmax) [3] improves L-Softmax by normalizing the weights, which achieves better performance on a series of open-set face recognition benchmarks [15, 16, 17]. Other loss functions [26, 27, 28, 29] based on contrastive loss or center loss also demonstrate effective performance on enhancing discrimination.

Normalization Approaches. Normalization has been studied in recent deep face recognition studies. [30] normalizes the weights which replace the inner product with cosine similarity within softmax loss. [31] applies the L2 constraint on features to embed faces in normalized space. Note that normalization on feature vectors or weight vec-

tors achieves much lower intra-class angular variability by concentrating more on the angle during training. Hence the angles between identities can be well optimized. The aforementioned A-Softmax [3] also adopts normalization approaches in feature learning.

3. The Proposed Approach

In the following of this section, we firstly introduce the proposed LMCL in detail (Sec. 3.1). And a comparison with other loss functions is given to show the superiority of LMCL over the other loss functions (Sec. 3.2). The feature normalization technique in LMCL is further described to clarify its effectiveness (Sec. 3.3). Lastly, we derive a theoretical analysis for the proposed LMCL. (Sec. 3.4).

3.1. Large Margin Cosine Loss

We start by rethinking the softmax loss from cosine perspective. The softmax loss separates features from different classes by maximizing posterior probability of the ground-truth class. Given an input feature vector x_i with its corresponding label y_i , the softmax loss can be formulated as:

$$L_s = \frac{1}{N} \sum_{i=1}^N -\log p_i = \frac{1}{N} \sum_{i=1}^N -\log \frac{e^{f_{y_i}}}{\sum_{j=1}^C e^{f_j}} \quad (1)$$

where p_i denotes the posterior probability of x_i being correctly classified. The N is the number of training samples and C is the number of classes. The f_j is usually formulated as activation of a fully connected layer with weight vector W_j and bias B_j . We fix the bias $B_j = 0$ for simplicity, and as a result f_j is given by:

$$f_j = W_j^T x = \|W_j\| \|x\| \cos \theta_j \quad (2)$$

where θ_j is the angle between W_j and x . This formula suggests that both norm and angle of vectors contribute to the posterior probability.

To develop effective feature learning, the norm of W should be necessarily invariable. We fix $\|W_j\| = 1$ by L2 normalization. At the testing stage, the face recognition score of a testing face pair is usually calculated according to cosine similarity between the two face features. This suggests the norm of feature x is not contributing to the scoring function. Thus, at the training stage, we fix $\|x\| = s$. Consequently, the posterior probability only relies on cosine of angle. The modified loss can be formulated as

$$L_{ns} = \frac{1}{N} \sum_i -\log \frac{e^{s \cos(\theta_{y_i, i})}}{\sum_j e^{s \cos(\theta_{j, i})}} \quad (3)$$

Because we remove variation in radial direction by fixing $\|x\| = s$, the resulting model learns features that are separable in angular space. We refer to this loss as the Normalized version of Softmax Loss (NSL) in this paper.

However, features learned by the NSL are not sufficiently discriminative because the NSL only emphasizes correctly classification. To address this, we introduce the cosine margin to the classification boundary by introducing the cosine margin, which is naturally incorporated with the cosine formulation of Softmax.

Considering a scenario of binary-classes for example, let θ_i denote the angle between the learned feature and the weight vector of Class C_i . The NSL forces $\cos(\theta_1) > \cos(\theta_2)$ for C_1 , and similarly for C_2 , so that features from different classes are correctly classified. To develop large margin classifier, we further require $\cos(\theta_1) - m > \cos(\theta_2)$ and $\cos(\theta_2) - m > \cos(\theta_1)$, where $m \geq 0$ is a fixed parameter introduced to control the magnitude of cosine margin. Since $\cos(\theta_i) - m$ is lower than $\cos(\theta_i)$, the constraint are more stringent for classification. The above analysis can be well generalized to the scenario of multi-classes. Therefore, the altered loss reinforces the discrimination of learned features by encouraging an extra margin in cosine space.

Formally, we define Large Margin Cosine Loss (LMCL) as:

$$L_{lmc} = \frac{1}{N} \sum_i -\log \frac{e^{s(\cos(\theta_{y_i, i}) - m)}}{e^{s(\cos(\theta_{y_i, i}) - m)} + \sum_{j \neq y_i} e^{s \cos(\theta_{j, i})}} \quad (4)$$

subject to

$$\begin{aligned} W &= \frac{W^*}{\|W^*\|} \\ x &= \frac{x^*}{\|x^*\|} \\ \cos(\theta_j, i) &= W_j^T x_i \end{aligned} \quad (5)$$

where the N is the number of training samples, the x_i is the i -th feature vector corresponding to the ground-truth class of y_i , the W_j is the weight vector of the j -th class, and the θ_j is the angle between W_j and x_i .

3.2. Comparison on Different Loss Functions

In this section, we compare the decision margin of our method (LMCL) to: Softmax, NSL, and A-Softmax, as illustrated in Figure 2. For simplicity of analysis, we consider the two classes case with class C_1 and C_2 . Let W_1 and W_2 denote weight vectors for C_1 and C_2 .

Softmax loss defines decision boundary by:

$$\|W_1\| \cos(\theta_1) = \|W_2\| \cos(\theta_2)$$

Thus, its boundary depends on both magnitude of weight vectors and cosine of angles, which results in overlapping decision area (margin < 0) in cosine space. This is illustrated in the first subplot of Figure 2. As noted before, at testing stage it is a common practice to only consider cosine similarity between testing features for face recognition.

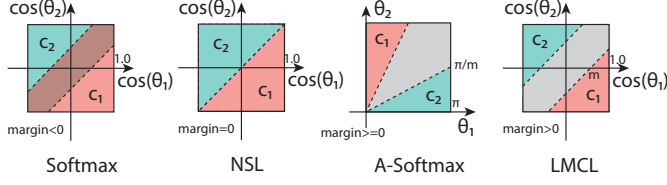


Figure 2. The comparison of decision margin for different loss functions for two classes. The dashed line represents decision boundary, and the gray areas are decision margins.

Consequently, the trained classifier with Softmax loss is unable to perfectly classify testing samples in cosine space.

NSL normalizes weight vectors W_1 and W_2 so that they have constant magnitude 1, which results in decision boundary given by:

$$\cos(\theta_1) = \cos(\theta_2)$$

The decision boundary of NSL is illustrated in the second subplot of Figure 2. We can see that by removing radial variation, the NSL is able to perfectly classify testing samples in cosine space, with margin = 0. However, it is not quite robust to noise because there is no decision margin: any small perturbation around decision boundary can change the decision result.

A-Softmax improves the softmax by introducing extra margin, such that the decision boundary is given by:

$$\begin{aligned} C_1 : \cos(m\theta_1) &\geq \cos(\theta_2) \\ C_2 : \cos(m\theta_2) &\geq \cos(\theta_1) \end{aligned}$$

Thus, for C_1 it requires $\theta_1 \leq \frac{\theta_2}{m}$, and similarly for C_2 . The third plot of Figure 2 depicts this decision area, where gray area denotes decision margin. However, the margin of A-Softmax is not consistent over all θ values: the margin becomes smaller as θ reduces, and it vanishes completely when $\theta = 0$. This results in two potential issues. First, for difficult classes C_1 and C_2 which are usually visually similar and thus smaller angle between W_1 and W_2 , the margin is consequently smaller. Second, technically speaking one has to employ extra facility with an ad-hoc piecewise function to overcome the non-monotonicity difficulty of the cosine function.

LMCL (proposed) defines decision margin in cosine space rather than angle space (like A-Softmax) by:

$$\begin{aligned} C_1 : \cos(\theta_1) &\geq \cos(\theta_2) + m \\ C_2 : \cos(\theta_2) &\geq \cos(\theta_1) + m \end{aligned}$$

Therefore, $\cos(\theta_1)$ is maximized while $\cos(\theta_2)$ being minimized for C_1 (similarly for C_2) to perform the large-margin classification. The last subplot in Figure 2 illustrates the decision boundary of LMCL in cosine space, we can see clear margin($\sqrt{2}m$) in the produced distribution of the cosine of

angles. This suggests that the LMCL is more robust than the NSL, because small perturbation around decision boundary (dashed line) less likely leads to incorrect decision. The cosine margin is applied consistently to all samples, regardless of the angles of their weight vectors.

3.3. Normalization on Features

In the proposed LMCL, normalization scheme is involved on purpose to derive the formulation of cosine loss and remove variation in radial direction. Unlike [3] that only normalizes the weight vectors, our approach simultaneously normalizes both the weight vectors and the feature vectors. As a result, the features distribute on a hypersphere where the scaling parameter s controls the magnitude of radius. In this section, we discuss why feature normalization is necessary and how feature normalization encourages better feature learning in the proposed LMCL approach.

The necessity of feature normalization is presented in two respects: First off, the original softmax loss without feature normalization implicitly learns both the Euclidean norm (L2-norm) of feature vectors and the cosine value of the angle. The L2-norm is adaptively learned for minimizing the overall loss, which results in the relatively weak cosine constraint. Particularly, the adaptive L2-norm of easy samples becomes much higher than hard samples to remedy the inferior performance of cosine metric. On the contrary, our approach requires the entire set of features to have the same Euclidean norm so that the learning process of features only depends on the cosine value to develop the discriminative power. Features from same classes are clustered together and those from different classes are pulled apart on surface of the hypersphere. Additionally, we consider the situation when the model initially starts to minimize the LMCL. Given the feature x , let $\cos(\theta_i)$, $\cos(\theta_j)$ denote cosine scores of the two classes respectively. Without normalization on features, the LMCL forces $\|x\|(\cos(\theta_i) - m) > \|x\| \cos(\theta_j)$. Note $\cos(\theta_i)$ and $\cos(\theta_j)$ can be initially comparable with each other. Thus, as long as $(\cos(\theta_i) - m)$ is smaller than $\cos(\theta_j)$, $\|x\|$ is required to decrease for minimizing the loss, which degenerates the optimization. Therefore, feature normalization is critical under the supervision of LMCL, especially when the networks are trained from scratch. Likewise, it is more favorable to fix the scaling parameter s instead of adaptively learning.

Furthermore, the scaling parameter s should be set to a properly large value to yield better-performing features with lower training loss. For NSL, the loss continuously goes down with higher s while too small s leads to the insufficient convergence even no convergence. For LMCL, we also need adequately large s to ensure sufficient hyperspace for feature learning with expected large margin.

In the following, we show the parameter s should have a lower bound to obtain expected classification performance.

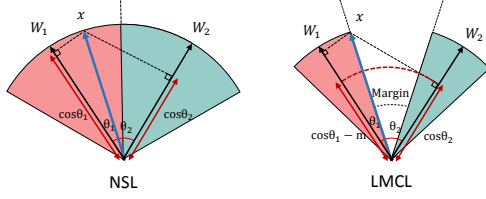


Figure 3. Geometrical interpretation of LMCL from feature perspective. Different color areas represent feature space from distinct classes. LMCL has a relatively compact feature region compared with NSL.

Given the normalized learned features x and unit weight vectors W , we denote the total number of classes as C . Suppose the learned features separately lie on the surface of hypersphere and center around the corresponding weight vector. Let P_w denotes the expected minimum posterior probability of class center (i.e. W), the lower bound of s is as follows²:

$$s \geq \frac{C-1}{C} \ln \frac{(C-1)P_w}{1-P_w}. \quad (6)$$

Based on this bound, we can infer that s should be enlarged consistently if we expect an optimal P_w for classification with a certain number of class. Besides, by keeping a fixed P_w , the desired s should be larger to deal with more classes since the growing number of classes increase the difficulty for classification in the relatively compact space. A hypersphere with large radius s is therefore required for embedding features with small intra-class distance and large inter-class distance.

3.4. Theoretical Analysis for LMCL

The preceding sections essentially discuss the LMCL on the classification point of view. In terms of learning the discriminative features on the hypersphere, cosine margin servers as momentous part to strengthen the discrimination power of features. A detailed analysis about the quantitative feasible choice of the cosine margin (i.e. the bound of hyper-parameter m) is necessary. The optimal choice of m potentially leads to more promising learning of highly discriminative face features. In the following section, we delve into the decision boundary and angular margin in the feature space to derive the theoretical bound for hyper-parameter m .

First off, considering the binary-classes case with class C_1 and C_2 as before, suppose the normalized feature vector x is given. Let W_i denotes the normalized weight vector, and θ_i denotes the angle between x and W_i . For NSL, the decision boundary defines as $\cos \theta_1 - \cos \theta_2 = 0$, which is equivalent to the angular bisector of W_1 and W_2 as shown

in the left of Figure 3. This addresses that the model supervised by NSL partitions the underlying feature space to two close regions where the features near the boundary are extremely ambiguous (i.e. belonging to either class is acceptable). In contrast, LMCL drives the decision boundary formulated by $\cos \theta_1 - \cos \theta_2 = m$ for C_1 , in which θ_1 should be much smaller than θ_2 (similarly for C_2). Consequently, the inter-class variance are enlarged while the intra-class variance is compacted.

Back to Figure 3, one can observe that the maximum angular margin is subject to the angle between W_1 and W_2 . Accordingly, the cosine margin should have the limited variable scope when the W_1 and W_2 are given. Specifically, imagine a scenario that all the feature vectors belonging to class i exactly overlap with the corresponding weight vector W_i of class i . In other words, every feature vector is identical to the weight vector for class i , and apparently the feature space is in an extreme situation where all the features lie on their class center. In that case, the margin of decision boundaries has been maximized (i.e. strict upper bound of cosine margin).

To extend in the general case, we suppose all the features are well-separated and we have a total number of C classes, the theoretical variable scope of m is supposed to be: $0 \leq m \leq (1 - \max(W_i^T W_j))$, where $i, j \leq n, i \neq j$. The softmax loss tries to maximize the angle between any of the two weight vectors from two different classes in order to perform perfect classification. Hence it is clear that the optimal solution for softmax loss should be to uniformly distribute the weight vectors on a unit hypersphere. Based on this assumption, the variable scope of introduced cosine margin m can be inferred as follows³:

$$\begin{aligned} 0 \leq m &\leq 1 - \cos \frac{2\pi}{C}, \quad (K = 2) \\ 0 \leq m &\leq \frac{C}{C-1}, \quad (C \leq K+1) \\ 0 \leq m &\ll \frac{C}{C-1}, \quad (C > K+1) \end{aligned} \quad (7)$$

where C is the number of training classes and K is the dimension of learned features. The inequalities indicate that as the number of classes increases, the upper bound of cosine margin between classes are correspondingly decreased. Especially, if the number of classes is much larger than feature dimension, the upper bound of cosine margin will be even smaller.

A reasonable choice of larger $m \in [0, \frac{C}{C-1})$ should effectively boost the learning of highly discriminative features. Nevertheless, the parameter m usually could not reach the theoretical upper bound in practice due to the vanish of feature space. That is, all the feature vectors are centered together according to the weight vector of the class.

²Proof is attached in supplemental material

³Proof is attached in supplemental material

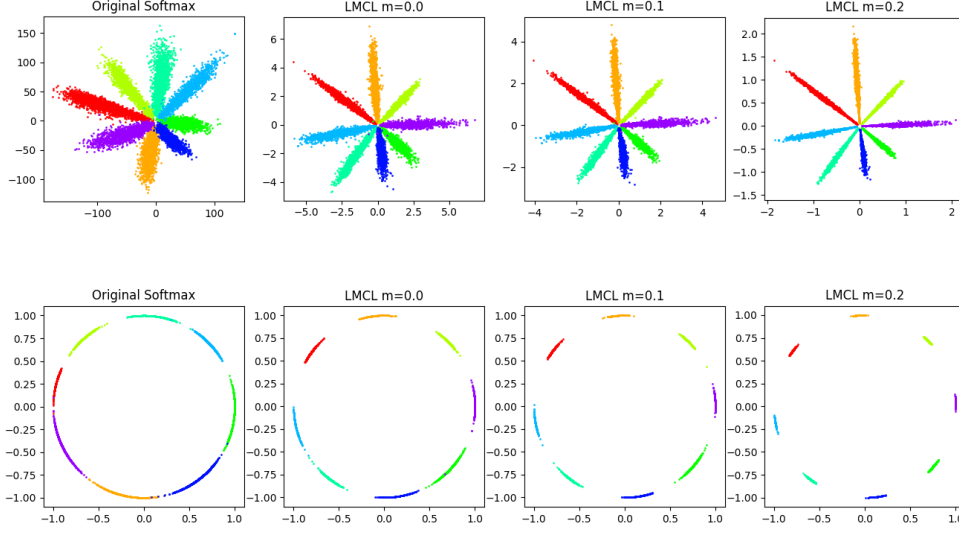


Figure 4. Toy experiment of different loss functions on 8 identities with 2D features. The first row maps the 2D features on Euclidean space while the second row projects the 2D features onto the angular space. The gap becomes evident as the margin term m increases.

In fact, the model fails to converge when m is too large, because the cosine constraint (i.e. $\cos \theta_1 - m > \cos \theta_2$ or $\cos \theta_2 - m > \cos \theta_1$ for two classes) becomes stricter and is hardly able to be satisfied. Besides, the cosine constraint with overlarge m forces the training process to be more sensitive to noise data. The ever-increasing m starts to degrade the overall performance at some point because of failing to converge.

We perform a toy experiment for better visualization on features and validating our algorithm. We select face images from 8 distinct identities containing enough samples to clearly show the feature points on the plot. Several models are trained using the original softmax loss and the proposed LMCL with different settings of m . We extract 2-D feature of face images for simplicity. As discussed above, the m should be no larger than $1 - \cos \frac{\pi}{4}$ (about 0.29), so we set up three settings of the m for comparison, which are $m = 0$, $m = 0.1$ and $m = 0.2$. As shown in Figure 4, the first row and second row present the feature distribution in Euclidean space and angular space respectively. We can observe that the original softmax loss produces ambiguity in decision boundaries while LMCL performs much better. As m increases, the angular margin between different classes has been correspondingly enlarged.

4. Experiments

4.1. Implementation Details

Preprocessing. Firstly, face area and landmarks are detected by MTCNN [32] for the entire set of training and testing images. Then, the 5 facial points (two eyes, nose

and mouth corners) are adopted to perform similarity transformation. After that we obtain the cropped faces which are then resized to be 112×96 . Following [1, 3], each pixel (in $[0, 255]$) in RGB images is normalized by subtracting 127.5 then dividing by 128.

Training. For fair comparison, the CNN architecture used in our work is similar to [3], which has 64 convolutional layers and is based on residual units [7]. We use Caffe [33] to implement modifications in loss layers and run the models. The CNN models are trained through SGD, with a batch size of 64 on 8 GPUs. The scaling parameter s in equation (4) is set to 64 empirically. For a direct and fair comparison to the existing results that use small training datasets (less than 0.5M images and 20K subjects) [17, 18], we train our models on a small training dataset, which is the publicly available CASIA-WebFace [34] dataset containing 0.49M face images from 10,575 subjects. We also use a large training dataset to evaluate the performance of our approach for benchmark comparison with the state-of-the-art results (using large training dataset) on the benchmark face dataset. The large training dataset that we use in this study is composed of several public datasets and a private face dataset, containing about 5M images from more than 90K identities. The training faces are horizontally flipped for data augmentation. In our experiments we remove face images belong to identities that appear in the testing datasets.

Testing. At testing stage, features of original image and the flipped image are concatenated together to compose the final face representation. The cosine distance of features is computed as the similarity score. Finally, face verification and identification are conducted by thresholding and

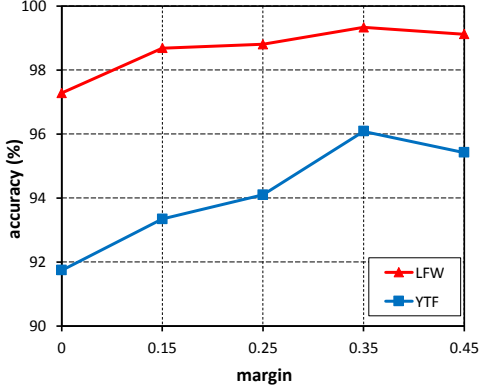


Figure 5. Performance (%) of CosFace with different margin parameters m on LFW[15] and YTF [16].

ranking the scores. We test our models on several popular public face datasets, including LFW[15], YTF[16], and MegaFace[17, 18].

4.2. Exploratory Experiments

Effect of m . The margin parameter m plays a key role in LMCL. In this part we conduct an experiment to investigate the effect of m . By varying m from 0 to 0.45 (If m is larger than 0.45, the model will fail to converge), we use the small training data (CASIA-WebFace [34]) to train our CosFace model and evaluate its performance on the LFW[15] and YTF[16] datasets, as illustrated in Figure 5. We can see that the model without the margin (in this case $m=0$) leads to the worst performance. As m being increased, the accuracies are improved consistently on both datasets, and get saturated at $m = 0.35$. This demonstrates the effectiveness of the margin m . By increasing the margin m , the discriminative power of the learned features can be significantly improved. In this study, m is set to fixed 0.35 in the subsequent experiments.

Effect of Feature Normalization. To investigate the effect of the feature normalization scheme in our approach, we train our CosFace models with and without the feature normalization scheme, and compare their performance on LFW[15], YTF[16], and the Megaface Challenge 1(MF1)[17]. Note that the model trained without normalization is initialized by softmax loss and then supervised by the proposed LMCL. The comparative results are reported in Table 1. It is very clear that the model using the feature normalization scheme consistently outperforms the model without the feature normalization scheme across the three datasets. As discussed above, feature normalization removes radical variance, and the learned features can be more discriminative in angular space. This experiment verifies this point.

Normalization	LFW	YTF	MF1 Rank 1	MF1 Veri.
No	99.10	93.1	75.10	88.65
Yes	99.33	96.1	77.11	89.88

Table 1. Comparison of our models with and without feature normalization on Megaface Challenge 1 (MF1). "Rank 1" refers to rank-1 face identification accuracy and "Veri." refers to face verification TAR (True Accepted Rate) under 10^{-6} FAR (False Accepted Rate).

Method	LFW	YTF	MF1 Rank1	MF1 Veri.
Softmax Loss[3]	97.88	93.1	54.85	65.92
Softmax+Contrastive[20]	98.78	93.5	65.21	78.86
Triplet Loss[23]	98.70	93.4	64.79	78.32
L-Softmax Loss[2]	99.10	94.0	67.12	80.42
Softmax+Center Loss[1]	99.05	94.4	65.49	80.14
A-Softmax[3]	99.42	95.0	72.72	85.56
A-Softmax-NormFea	99.32	95.4	75.42	88.82
LMCL	99.33	96.1	77.11	89.88

Table 2. Comparison of the proposed LMCL with the state-of-the-art loss functions in face recognition community. All the methods in this table are using the same training data and the same 64-layer CNN architecture.

4.3. Comparison with the state-of-the-art loss functions

In this part, we compare the performance of the proposed LMCL with the state-of-the-art loss functions in face recognition community. Following the experimental setting in [3], we train a model with the guidance of the proposed LMCL on the CAISA-WebFace[34] using the same 64-layer CNN architecture described in [3]. The experimental comparison on LFW, YTF and MF1 are reported in Table 2. For fair comparison, we are strictly following the model structure (a 64-layers ResNet-Like CNNs) and the detailed experimental settings of SphereFace [3]. As can be seen in Table 2, LMCL consistently achieves competitive results compared to the other losses across the three datasets. Especially, our method not only surpasses the performance of A-Softmax with feature normalization (named as A-Softmax-NormFea in Table 2), but also significantly outperforms the other loss functions on YTF and MF1, which demonstrates the effectiveness of LMCL.

4.4. Overall Benchmark Comparisons

4.4.1 Evaluation on LFW and YTF

LFW [15] is a standard face verification testing dataset in unconstrained conditions. It includes 13,233 face images from 5749 identities collected from the website. We evaluate our model strictly following the standard protocol of unrestricted with labeled outside data [15], and report the result on the 6,000 pair testing images. YTF [16] contains 3,425 videos of 1,595 different people. The average

Method	Training Data	#Models	LFW	YTF
Deep Face[9]	4M	3	97.35	91.4
FaceNet[23]	200M	1	99.63	95.1
DeepFR [35]	2.6M	1	98.95	97.3
DeepID2+[21]	300K	25	99.47	93.2
Center Face[1]	0.7M	1	99.28	94.9
Baidu[36]	1.3M	1	99.13	-
SphereFace[3]	0.49M	1	99.42	95.0
CosFace	5M	1	99.73	97.6

Table 3. Face verification (%) on the LFW and YTF datasets. “#Models” indicates the number of models have been used in the method for evaluation.

Method	Protocol	MF1 Rank1	MF1 Veri.
SIAT_MMLAB[1]	Small	65.23	76.72
DeepSense - Small	Small	70.98	82.85
SphereFace - Small[3]	Small	75.76	90.04
Beijing FaceAll V2	Small	76.66	77.60
GRCCV	Small	77.67	74.88
FUDAN-CS_SDS[24]	Small	77.98	79.19
CosFace(Single-patch)	Small	77.11	89.88
CosFace(3-patch ensemble)	Small	79.54	92.22
Beijing FaceAll_Norm_1600	Large	64.80	67.11
Google - FaceNet v8[23]	Large	70.49	86.47
NTechLAB - facenx_large	Large	73.30	85.08
SIATMMLAB TencentVision	Large	74.20	87.27
DeepSense V2	Large	81.29	95.99
YouTu Lab	Large	83.29	91.34
Vocord - deepVo V3	Large	91.76	94.96
CosFace(Single-patch)	Large	82.72	96.65
CosFace(3-patch ensemble)	Large	84.26	97.96

Table 4. Face Identification and Verification Evaluation on Megaface 2 (MF2). “Rank 1” refers to rank-1 face identification accuracy and “Veri.” refers to face verification TAR (True Accepted Rate) under 10^{-6} FAR (False Accepted Rate).

Method	Protocol	MF2 Rank1	MF2 Veri.
3DiVi	Large	57.04	66.45
Team 2009	Large	58.93	71.12
NEC	Large	62.12	66.84
GRCCV	Large	75.77	74.84
SphereFace	Large	71.17	84.22
CosFace (Single-patch)	Large	74.11	86.77
CosFace(3-patch ensemble)	Large	77.06	90.30

Table 5. Face Identification and Verification Evaluation on Megaface Challenge 2 (MF2). “Rank 1” refers to rank-1 face identification accuracy and “Veri.” refers to face verification TAR (True Accepted Rate) under 10^{-6} FAR (False Accepted Rate).

length of a video clip is 181.3 frames. All the video sequences were downloaded from YouTube. We follow the unrestricted with labeled outside data protocol and report the result on 5,000 video pairs.

As shown in Table 3, the proposed CosFace achieves the state-of-the-art results of 99.73% on LFW and 97.6% on YTF. FaceNet achieves the runner-up performance on LFW with the large scale of the image dataset, which has approxi-

mately 200 million face images. In terms of YTF, our model reaches the first place over all other methods.

4.4.2 Evaluation on Megaface

MegaFace [17, 18] is a very challenging testing benchmark recently released for large-scale face identification and verification, which contains a gallery set and a probe set. The gallery set in Megaface is composed of more than 1 million face images. The probe set has two existing databases: Facescrub [37] and FGNET [38]. The facescrub dataset contains 106,863 face images of 530 celebrities. The FGNET dataset is a relatively small dataset (including 1002 face images from 82 persons with each one having multiple face images at different ages), which is mainly used for testing age invariant face recognition. In this study, we use the Facescrub dataset as the probe set to evaluate the performance of our approach on both Megaface Challenge 1 and Challenge 2.

MegaFace Challenge 1 (MF1). On the MegaFace Challenge 1 [17], The gallery set incorporates more than 1 million images from 690K individuals collected from Flickr photos [39]. Table 4 summarizes the results of our models trained on two protocols of MegaFace where the training dataset is regarded as small if it has less than 0.5 million images, large otherwise. The CosFace approach shows its superiority for both the identification and verification tasks on both the protocols.

MegaFace Challenge 2 (MF2). In terms of MegaFace Challenge 2 [18], all the algorithms need to use the training data provided by MegaFace. The training data for Megaface Challenge 2 contains 4.7 million faces and 672K identities, which corresponds to the large protocol. The gallery set has 1 million images that are different from the challenge 1 gallery set. Not surprisingly, Our method wins the first place of challenge 2 in table 5, setting a new state-of-the-art with a large margin (1.39% on rank-1 identification accuracy and 5.46% on verification performance).

5. Conclusion

In this paper, we propose an innovative approach named LMCL to guide the deep CNNs to learn highly discriminative features by extending the cosine margin between decision boundaries, for boosting the performance of deep face recognition. We provide the well-formed geometrical and theoretical interpretation to verify the effectiveness of the proposed LMCL on generating decent face representation, following by extensive experiments on various datasets. Our approach consistently achieves the state-of-the-art results across several face benchmarks as depicted in experiment section. We wish the substantial explorations on face feature will benefit the face recognition community and LMCL can be a valuable pioneering work in the future.

References

- [1] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. A discriminative feature learning approach for deep face recognition. In *European Conference on Computer Vision (ECCV)*, pages 499–515, 2016. 1, 2, 6, 7, 8
- [2] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. Large-Margin Softmax Loss for Convolutional Neural Networks. In *International Conference on Machine Learning (ICML)*, 2016. 1, 2, 7
- [3] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. SphereFace: Deep Hypersphere Embedding for Face Recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 2, 3, 4, 6, 7, 8
- [4] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations (ICLR)*, 2015. 1, 2
- [5] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2012. 1, 2
- [6] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-Excitation Networks. In *arXiv preprint arXiv:1709.01507*, 2017. 1
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 6
- [8] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. *arXiv preprint arXiv:1611.05431*, 2016. 1
- [9] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 1, 2, 8
- [10] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation from predicting 10,000 classes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 1, 2
- [11] Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively, with application to face verification. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005. 1, 2
- [12] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2006. 1, 2
- [13] Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *International Workshop on Similarity-Based Pattern Recognition*, 2015. 1, 2
- [14] Jiang Wang, Yang Song, Thomas Leung, Chuck Rosenberg, Jingbin Wang, James Philbin, Bo Chen, and Ying Wu. Learning fine-grained image similarity with deep ranking. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 1, 2
- [15] Gary B Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In *Technical Report 07-49, University of Massachusetts, Amherst*, 2007. 2, 6, 7
- [16] Lior Wolf, Tal Hassner, and Itay Maoz. Face recognition in unconstrained videos with matched background similarity. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011. 2, 6, 7
- [17] Ira Kemelmacher-Shlizerman, Steven M Seitz, Daniel Miller, and Evan Brossard. The megaface benchmark: 1 million faces for recognition at scale. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 6, 7, 8
- [18] Aaron Nech and Ira Kemelmacher-Shlizerman. Level playing field for million scale face recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2, 6, 7, 8
- [19] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2
- [20] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In *Advances in Neural Information Processing Systems (NIPS)*, 2014. 2, 7
- [21] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deeply learned face representations are sparse, selective, and robust. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2, 8
- [22] Yi Sun, Ding Liang, Xiaogang Wang, and Xiaoou Tang. Deepid3: Face recognition with very deep neural networks. In *arXiv preprint arXiv:1502.00873*, 2015. 2
- [23] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2, 7, 8
- [24] Zhanxiong Wang, Keke He, Yanwei Fu, Rui Feng, Yu-Gang Jiang, and Xiangyang Xue. Multi-task Deep Neural Network for Joint Face Recognition and Facial Attribute Prediction. In *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval (ICMR)*, 2017. 2, 8
- [25] Guosheng Hu, Yongxin Yang, Dong Yi, Josef Kittler, William Christmas, Stan Z Li, and Timothy Hospedales. When face recognition meets with deep learning: an evaluation of convolutional neural networks for face recognition. In *International Conference on Computer Vision Workshops (ICCVW)*, 2015. 2
- [26] Xiao Zhang, Zhiyuan Fang, Yandong Wen, Zhifeng Li, and Yu Qiao. Range Loss for Deep Face Recognition with Long-tail. In *International Conference on Computer Vision (ICCV)*, 2017. 2

- [27] Jiankang Deng, Yuxiang Zhou, and Stefanos Zafeiriou. Marginal loss for deep face recognition. In *Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2017. 2
- [28] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: a deep quadruplet network for person re-identification. In *arXiv preprint arXiv:1704.01719*, 2017. 2
- [29] Jie Cai, Zibo Meng, Ahmed Shehab Khan, Ziyuan Li, and Yan Tong. Island Loss for Learning Discriminative Features in Facial Expression Recognition. In *arXiv preprint arXiv:1710.03144*, 2017. 2
- [30] Feng Wang, Xiang Xiang, Jian Cheng, and Alan L Yuille. NormFace: L_2 Hypersphere Embedding for Face Verification. In *Proceedings of the 2017 ACM on Multimedia Conference (ACM MM)*, 2017. 2
- [31] Rajeev Ranjan, Carlos D Castillo, and Rama Chellappa. L2-constrained Softmax Loss for Discriminative Face Verification. In *arXiv preprint arXiv:1703.09507*, 2017. 2
- [32] K. Zhang, Z. Zhang, Z. Li and Y. Qiao. Joint Face Detection and Alignment using Multi-task Cascaded Convolutional Networks. *Signal Processing Letters*, 23(10):1499–1503, 2016. 6
- [33] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 2016 ACM on Multimedia Conference (ACM MM)*, 2014. 6
- [34] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. In *arXiv preprint arXiv:1411.7923*, 2014. 6, 7
- [35] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, et al. Deep face recognition. In *BMVC*, volume 1, page 6, 2015. 8
- [36] Jingtuo Liu, Yafeng Deng, Tao Bai, Zhengping Wei, and Chang Huang. Targeting ultimate accuracy: Face recognition via deep embedding. *arXiv preprint arXiv:1506.07310*, 2015. 8
- [37] Hong-Wei Ng and Stefan Winkler. A data-driven approach to cleaning large face datasets. In *Image Processing (ICIP), 2014 IEEE International Conference on*, pages 343–347. IEEE, 2014. 8
- [38] *FG-NET Aging Database*, <http://www.fgnet.rsunit.com/>. 8
- [39] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. YFCC100M: The new data in multimedia research. *Communications of the ACM*, 2016. 8

Appendix: Supplementary Material

This supplementary document provides mathematical details for the derivation of the lower bound of the scaling parameter s (Equation 6 in the main paper), and the variable scope of the cosine margin m (Equation 7 in the main paper).

Proposition of the Scaling Parameter s

Given the normalized learned features x and unit weight vectors W , we denote the total number of classes as C where $C > 1$. Suppose that the learned features separately lie on the surface of a hypersphere and center around the corresponding weight vector. Let P_w denote the expected minimum posterior probability of the class center (*i.e.*, W). The lower bound of s is formulated as follows:

$$s \geq \frac{C-1}{C} \ln \frac{(C-1)P_w}{1-P_w}$$

Proof:

Let W_i denote the i -th unit weight vector. $\forall i$, we have:

$$\frac{e^s}{e^s + \sum_{j,j \neq i} e^{s(W_i^T W_j)}} \geq P_w, \quad (8)$$

$$1 + e^{-s} \sum_{j,j \neq i} e^{s(W_i^T W_j)} \leq \frac{1}{P_w}, \quad (9)$$

$$\sum_{i=1}^C (1 + e^{-s} \sum_{j,j \neq i} e^{s(W_i^T W_j)}) \leq \frac{C}{P_w}, \quad (10)$$

$$1 + \frac{e^{-s}}{C} \sum_{i,j,i \neq j} e^{s(W_i^T W_j)} \leq \frac{1}{P_w}. \quad (11)$$

Because $f(x) = e^{s \cdot x}$ is a convex function, according to Jensen's inequality, we obtain:

$$\frac{1}{C(C-1)} \sum_{i,j,i \neq j} e^{s(W_i^T W_j)} \geq e^{\frac{s}{C(C-1)} \sum_{i,j,i \neq j} W_i^T W_j}. \quad (12)$$

Besides, it is known that

$$\sum_{i,j,i \neq j} W_i^T W_j = (\sum_i W_i)^2 - (\sum_i W_i^2) \geq -C. \quad (13)$$

Thus, we have:

$$1 + (C-1)e^{-\frac{sC}{C-1}} \leq \frac{1}{P_w}. \quad (14)$$

Further simplification yields:

$$s \geq \frac{C-1}{C} \ln \frac{(C-1)P_w}{1-P_w}. \quad (15)$$

The equality holds if and only if every $W_i^T W_j$ is equal ($i \neq j$), and $\sum_i W_i = 0$. Because at most $K+1$ unit vectors are able to satisfy this condition in the K -dimension hyper-space, the equality holds only when $C \leq K+1$, where K is the dimension of the learned features.

Proposition of the Cosine Margin m

Suppose that the weight vectors are uniformly distributed on a unit hypersphere. The variable scope of the introduced cosine margin m is formulated as follows :

$$0 \leq m \leq 1 - \cos \frac{2\pi}{C}, \quad (K=2)$$

$$0 \leq m \leq \frac{C}{C-1}, \quad (K>2, C \leq K+1)$$

$$0 \leq m \ll \frac{C}{C-1}, \quad (K>2, C > K+1)$$

where C is the total number of training classes and K is the dimension of the learned features.

Proof:

For $K = 2$, the weight vectors uniformly spread on a unit circle. Hence, $\max(W_i^T W_j) = \cos \frac{2\pi}{C}$. It follows $0 \leq m \leq (1 - \max(W_i^T W_j)) = 1 - \cos \frac{2\pi}{C}$.

For $K > 2$, the inequality below holds:

$$C(C-1) \max(W_i^T W_j) \geq \sum_{i,j,i \neq j} W_i^T W_j = (\sum_i W_i)^2 - (\sum_i W_i^2) \geq -C. \quad (16)$$

Therefore, $\max(W_i^T W_j) \geq \frac{-1}{C-1}$, and we have $0 \leq m \leq (1 - \max(W_i^T W_j)) \leq \frac{C}{C-1}$.

Similarly, the equality holds if and only if every $W_i^T W_j$ is equal ($i \neq j$), and $\sum_i W_i = 0$. As discussed above, this is satisfied only if $C \leq K + 1$. On this condition, the distance between the vertexes of two arbitrary W should be the same. In other words, they form a regular simplex such as an equilateral triangle if $C = 3$, or a regular tetrahedron if $C = 4$.

For the case of $C > K + 1$, the equality cannot be satisfied. In fact, it is unable to formulate the strict upper bound. Hence, we obtain $0 \leq m \ll \frac{C}{C-1}$. Because the number of classes can be much larger than the feature dimension, the equality cannot hold in practice.