

# SphereFace and CosFace

Zhuonan Lin, A53218121

May 30, 2019

## A SphereFace for Face Verification

### A.1 Evaluation SphereFace on LFW dataset

As the homework instructs, we download the LFW [1] dataset and the PyTorch code for ShpereFace [2], provided on the course website. Unzipping both the dataset and pretrained model, we can evaluate the SphereFace model on the LFW dataset by running:

```
python ./lfw_eval.py --model model/sphere20a_20171020.pth
```

The evaluation process and final result will be printed on the screen. The final accuracy for pretrained SphereFace on LFW dataset is 0.9918, with standard derivation and threshold of different folds are 0.0051 and 0.3095, respectively.

### A.2 Evaluation Details on LFW dataset

#### A.2.1 Distance Metric

In *lfw\_eval.py*: Line 135, the distance metric to measure the distance between two faces is calculated by:

```
1 cosdistance = f1.dot(f2)/(f1.norm()*f2.norm()+1e-5)
```

Here,  $f1$  and  $f2$  are the vectors for the two face images to be measured. This line of the code indicates that, as the variable name in the left hand side suggests, the distance between the two faces is measured by the cosine similarity, i.e. the cosine value of the angle between two face vectors. The small  $\epsilon = 1e - 5$  in the denominator is set to prevent dividing by zero. This metric is corresponding to [2], as it introduces in the **Testing** part in Section 4.1 in the paper:

”For all experiments, the final representation of a testing face is obtained by concatenating its original face features and its horizontally flipped features. The score (metric) is computed by the *cosine distance* of two features.”

## A.2.2 Threshold and Accuracy Computation

The threshold for verification determination and the test accuracy computation are done in *lfw\_eval.py*: Line 141-148

```
1 folds = KFold(n=pairNum, n_folds=10, shuffle=False)
2 thresholds = np.arange(-1.0, 1.0, 0.005)
3 predicts = np.array([*map(lambda line:line.strip('\n').split(),
    ... predicts) ] )
4 for idx, (train, test) in enumerate(folds):
5     best_thresh = find_best_threshold(thresholds, predicts[train]
    ... )
6     accuracy.append(eval_acc(best_thresh, predicts[test]))
7     thd.append(best_thresh)
8 print('LFWACC={:.4f} std={:.4f} thd={:.4f}'.format(np.mean(
    ... accuracy), np.std(accuracy), np.mean(thd)))
```

From the code snippet, we know that the evaluation is carried out with a K-Fold style. There is a threshold candidate pool ranges from -1.0 to 1.0, with a step 0.005. For each fold, the *find\_best\_threshold()* function will be called, where each threshold candidate will be enumerated and the corresponding accuracy on the training set will be evaluated. The threshold that yields the best accuracy will be returned and subsequently used as threshold on this fold.

Accuracy evaluation need to be performed both in threshold determination and final test set evaluation. Here, the accuracy for a given data threshold is calculated by the function *eval\_acc()*, where the predicted label with distance larger than the threshold will be regarded as 1 and otherwise 0. Finally, the accuracy will be calculated between the converted label and the ground truth and returned as the final accuracy result.

## A.3 Face Alignment and Crop

### A.3.1 Default Alignment and Crop

The face alignment in our code, is done by the function *alignment()* in *lfw\_eval.py*:Line 11-26

```
1 def alignment(src_img, src_pts):
2     ref_pts = [ [30.2946, 51.6963],[65.5318, 51.5014],
3                 [48.0252, 71.7366],[33.5493, 92.3655],[62.7299, 92.2041]
4                 ... ]
5     crop_size = (96, 112)
6     src_pts = np.array(src_pts).reshape(5,2)
7     s = np.array(src_pts).astype(np.float32)
8     r = np.array(ref_pts).astype(np.float32)
9
10    tfm = get_similarity_transform_for_PIL(s, r)
```

```

11
12     face_img = src_img.transform(crop_size, Image.AFFINE,
13                                 tfm.reshape(6), resample=Image.BILINEAR)
14     face_img = np.asarray(face_img)[: , : , ::-1]
15
16     return face_img

```

Here, the *ref\_pts* are the target positions of landmarks after alignment, which have been determined during the SphereFace training process. Then, the transform will be calculated using source points (*src\_pts*) and target points by the function *get\_similarity\_transform\_for\_PIL()*. Given this transform (*tmf*), the original image *src\_img* will be transformed, together with the bilinear interpolation and cropped into desired size (*crop\_size*). The final resultant image (*face\_img*) will be warped and returned.

### A.3.2 Crop Without Alignment

In this part, we will try to evaluate the SphereFace without face alignment. Instead, we will crop a image patch, with same size used in face alignment, at the center of each face image. In our code, we will first fill in the *cropping()* function defined in *lfw\_eval.py*: Line 28 as following, and two cropping examples for the alphabetically first two face images in LFW dataset are shown in Figure 1.

```

1  def cropping(src_img ):
2      # IMPLEMENT cropping the center of image
3      # Define crop_size as required
4      crop_size = (112, 96)
5      # convert original input to array
6      s = np.asarray(src_img)
7      # dicide image size (H: height and W: width)
8      H, W, _ = np.asarray(s).shape
9      # get the four boundaries index, center at image center
10     bottom, top = int(H/2.+crop_size[0]/2.), int(H/2.-crop_size
11         ... [0]/2.)
12     right, left = int(W/2.+crop_size[1]/2.), int(W/2.-crop_size
13         ... [1]/2.)
14     # crop with boudaries above, and convert to array for
15     ... returning
16     face_img = np.asarray(s[top:bottom, left:right, :])
17
18     return face_img

```

With this cropping function, we can then evaluate the ShpereFace on LFW dataset again without face alignment but center cropping by running:

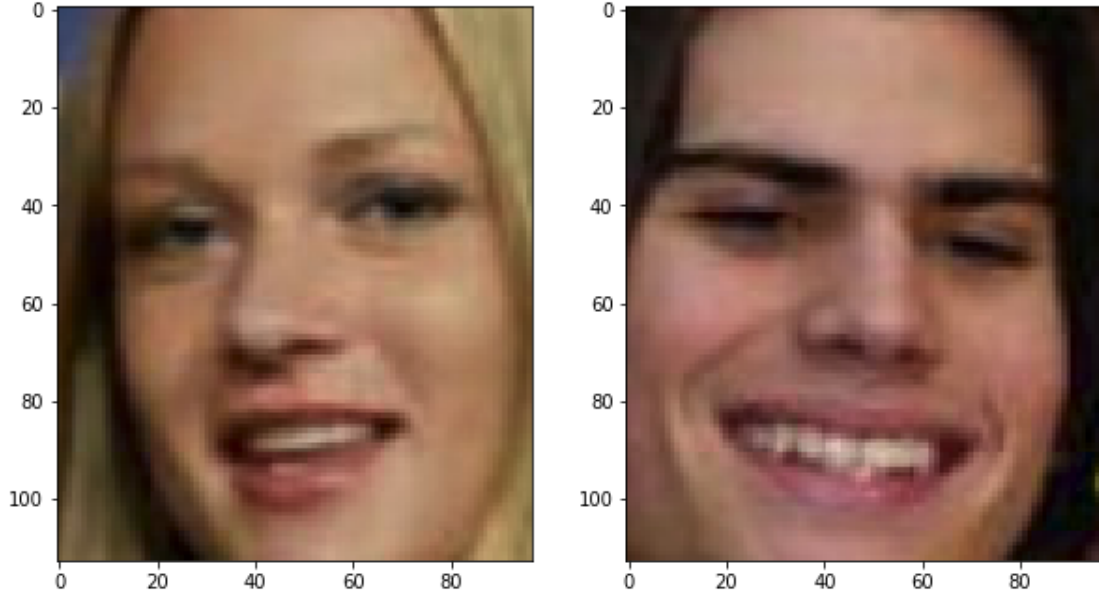


Figure 1: Cropping examples with center at image center and a height and width of 112 and 96 pixels, respectively.

```
python ./lfw_eval.py --model model/sphere20a_20171020.pth --alignmentMode 0
```

The evaluation process and final result will be printed on the screen. With this center cropping, the final accuracy for pretrained SphereFace on LFW dataset is 0.8297, with standard derivation and mean threshold of different folds are 0.0252 and 0.2440, respectively. The results are worse than those with face alignment as we can expect.

## B MTCNN Face Landmark Detection

### B.1 MTCNN Face Landmark Detection Examples

In this part, we use the MTCNN method [3] for pre-processing the face detection and alignment. After downloading the modified codes from course website, we run *lfw\_landmark.py* to perform MTCNN on LFW dataset. The results will be written in *../sphereFace/data/lfw\_landmarkMTCNN.txt*, in the form for further SphereFace inputting. Two example output lines are as following:

```
AJ_Cook/AJ_Cook_0001.jpg 103.672 114.958 146.303 109.461 118.223 135.837
110.155 158.646 149.281 154.148
AJ_Lamas/AJ_Lamas_0001.jpg 104.424 111.958 149.001 118.325 130.224 134.433
97.063 152.340 147.845 158.792
```

where the numbers following the image name are the coordinates of landmarks found by MTCNN. The detection results of the above two face images are shown below (this can be

done separately by a simple piece of codes using the *detect\_faces* and *show\_bboxes* methods from MTCNN).

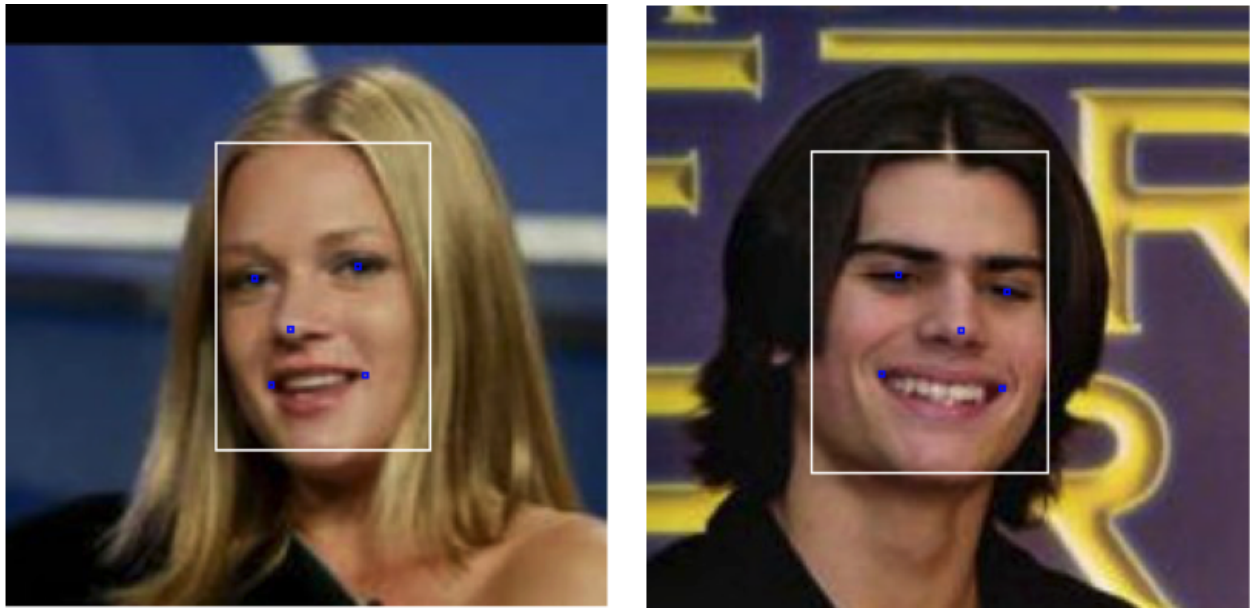


Figure 2: Face landmarks detection result examples using MTCNN. Landmarks are shown together with the bounding box also from MTCNN dectector.

## B.2 SphereFace with MTCNN Landmarks

With landmarks detected from MTCNN in the previous section, we can re-run the evaluation of SphereFace on LFW dataset, with the following line:

```
python ./lfw_eval.py --model model/sphere20a_20171020.pth --alignmentMode 2
```

The evaluation process and final result will be again printed on the screen. This time, the final accuracy with MTCNN pre-detected landmarks is 0.9845, with standard derivation and threshold of different folds are 0.0057 and 0.2995, respectively.

## B.3 Comparison to Original Results

Compare to the evaluation results to the one with default given face landmarks, our result with our MTCNN detected landmarks is not better, and even worse by 0.73%. One possible reason is that, although [2] claims that they also used MTCNN detected landmarks in the whole network, the details about the application is unknown to us. Their MTCNN model may be trained on different dataset so that the detected results are possibly slightly different from and more accurately than the pre-trained model we used, so that it out-performs a bit in A than B. Also, the crop procedure in the given pre-trained SphereFace is directly related

to the training process (e.g. the pre-determined *ref\_pts* in *alignment()* function). Therefore, we think there are some possible ways to improve our MTCNN + SphereFace results. First, we could retrain the MTCNN on a larger dataset for better performance. Second, we can try to relate MTCNN with SphereFace in the training stage by ourselves, which can give better reference points for cropping and alignment.

## B.4 Real-time Speed in MTCNN

One advantage of MTCNN is that the network can achieve the state-of-art accuracy in real-time. The real-time speed is achieved due to three main reasons: **carefully designed cascaded CNNs architecture, online hard sample mining strategy, and joint face alignment learning.**

First, The MTCNN completes its joint detection and alignment with a multistage cascaded CNN. There are 3 stages, where each has a individual CNN to achieve the goal step by step, namely, proposal network (P-Net, a fully convolutional network), refine network (R-Net) and output network (O-Net). P-Net will obtain all the candidate facial windows and their bounding box regression vectors. R-Net will refine the candidates and reject a large number of false ones. O-Net will identify face regions and five facial landmarks' positions. Note that the input for the whole structure is the image pyramid, which consists of different scales of the resized original image. In these three CNNs,  $3 \times 3$  **filters** are used instead of previous  $5 \times 5$  filters, which can reduce computing for less run time, while increase the depth to get better performance. Consequently, given this cascade structure with lightweight CNNs, MTCNN method can achieve real time speed in joint face detection and alignment task.

Second, in the MTCNN method, the authors came up with a online hard sample mining strategy in face/nonface classification task. To do this, in each minibatch, the losses computed in the forward propagation from all samples are sorted and the top 70% of them will be selected as hard samples. Only the gradients from these hard samples are calculated in the backward propagation. This online hard mining strategy ignores the easy samples that are less helpful to strengthen the detector during training and also help the runtime performance.

Last, the face alignment is learning jointly with detection. In training process, the three loss functions from individual CNN (face classification loss, bounding box regression loss, and facial landmark localization loss) will compose the final loss together as a multisource training. The final loss is tuned by some task-importance coefficients at different stages. This joint training strategy will simultaneously learn the detection and alignment tasks, and therefore improve the performance for less time.

## B.5 Non-maximal suppression (NMS) in MTCNN

Non-maximal suppression (NMS) is used both in P-Net and R-Net. In P-Net, NMS is employed non-maximum to merge highly overlapped face window candidates. In R-Net, NMS is used to help further reject a large number of false candidates. In the MTCNN code, the NMS procedures are carried out in *src/box\_utils.py*: Line 5-68.

```

1 def nms(boxes, overlap_threshold=0.5, mode='union'):
2     """Non-maximum suppression.
3
4     Arguments:
5         boxes: a float numpy array of shape [n, 5],
6               where each row is (xmin, ymin, xmax, ymax, score).
7         overlap_threshold: a float number.
8         mode: 'union' or 'min'.
9
10    Returns:
11        list with indices of the selected boxes
12    """
13
14    # if there are no boxes, return the empty list
15    if len(boxes) == 0:
16        return []
17
18    # list of picked indices
19    pick = []
20
21    # grab the coordinates of the bounding boxes
22    x1, y1, x2, y2, score = [boxes[:, i] for i in range(5)]
23
24    area = (x2 - x1 + 1.0)*(y2 - y1 + 1.0)
25    ids = np.argsort(score) # in increasing order
26
27    while len(ids) > 0:
28
29        # grab index of the largest value
30        last = len(ids) - 1
31        i = ids[last]
32        pick.append(i)
33
34        # compute intersections
35        # of the box with the largest score
36        # with the rest of boxes
37
38        # left top corner of intersection boxes
39        ix1 = np.maximum(x1[i], x1[ids[:last]])
40        iy1 = np.maximum(y1[i], y1[ids[:last]])
41
42        # right bottom corner of intersection boxes
43        ix2 = np.minimum(x2[i], x2[ids[:last]])

```

```

44         iy2 = np.minimum(y2[i], y2[ids[:last]])
45
46         # width and height of intersection boxes
47         w = np.maximum(0.0, ix2 - ix1 + 1.0)
48         h = np.maximum(0.0, iy2 - iy1 + 1.0)
49
50         # intersections' areas
51         inter = w * h
52         if mode == 'min':
53             overlap = inter/np.minimum(area[i], area[ids[:last]])
54         elif mode == 'union':
55             # intersection over union (IoU)
56             overlap = inter/(area[i] + area[ids[:last]] - inter)
57
58         # delete all boxes where overlap is too big
59         ids = np.delete(
60             ids,
61             np.concatenate([[last], np.where(overlap >
62                 ... overlap_threshold)[0]])
63         )
64     return pick

```

The basic idea for NMS used here is to **iteratively pick the current maximum score box** ( $boxes[i, :]$ ) and **eliminate the boxes whose overlap with  $boxes[i, :]$  is larger than the threshold** ( $overlap\_threshold$ ), until all candidates (in  $boxes$ ) are either picked or eliminated.

To be specific, in the code:

1. First we get the index list that has the order sorted by the score of each box ( $ids$ ). While we still have box that neither picked nor eliminated ( $len(ids) > 0$ ), we picked the current largest score box ( $ids[last]$ ) to be one of the remaining box to return (i.e.  $pick.append(ids[last])$ ).
2. Then we carry on the NMS filter based on the current largest score box. Given the current max score box, we find the index of left top and right bottom corner ( $ix1, iy1, ix2, iy2$ ) of the intersection area by calculating the maximum and minimum of  $x$  and  $y$  indexes among the highest score box and others. (Line: 38-48 above). Given  $ix1, iy1, ix2, iy2$ , **the ratio of intersection area to the box area** can be calculated. In the actual code, there are two types of intersection area, namely *min* and *union*, which calculates minimum overlap or the total union overlap, as the names suggest, respectively. Using the overlap area, each other box will be filter by whether its overlap with highest score box is larger than the threshold ( $overlap\_threshold$ ).
3. Repeat 1 and 2 until all boxes are either picked or eliminated. The picked boxes' indexes in  $pick$  will be returned, others are filtered by NMS method. This can reduced



the number of candidate facial window effectively, and thus help accelerate the overall speed of MTCNN.

## C Training SphereFace on CASIA Dataset

In this part, we will train the SphereFace using the angular softmax loss on CASIA dataset [4], and evaluate the performance on LFW dataset.

### C.1 A-Softmax Loss and $\psi$ Function

To train the neural network, we will follow the A-softmax loss function defined by Eq.(7) in [2]:

$$L_{ang} = \frac{1}{N} \sum_i -\log\left(\frac{e^{\|\mathbf{x}_i\| \psi(\theta_{y_i,i})}}{e^{\|\mathbf{x}_i\| \psi(\theta_{y_i,i})} + \sum_{j \neq y_i} e^{\|\mathbf{x}_i\| \cos(\theta_{j,i})}}\right) \quad (1)$$

where

$$\psi(\theta_{y_i,i}) = (-1)^k \cos(m\theta_{y_i,i}) - 2k \quad (2)$$

with  $\theta_{y_i,i} \in [\frac{k\pi}{m}, \frac{(k+1)\pi}{m}]$ ,  $k \in [0, m-1]$  and  $m \geq 1$ .

In *sphereface/faceNet.py*, we implement Eq.(2) under class *CustomLinear* and Eq.(1) under class *CustomLoss* respectively. The following snippets showing the the specific implementation with PyTorch:

*Implement psi\_theta*

```

1      # def forward(self, input):
2          # ...
3
4          # IMPLEMENT phi_theta
5          # Use multi-angle formulae to calculate cos(m*theta)
6          # these formulae given in self.mlambda
7          cos_m_theta = self.mlambda[self.m](cos_theta)
8          # Calculate the hyper-parameter k
9          # k should be in the [0, m-1] and
10         # theta should be in [k*pi/m, (k+1)*pi/m]
11         # here in the forward process, we choose k that makes
12         ... theta = k*pi/m
13         theta = Variable(cos_theta.data.acos())
14         k = (self.m * theta / np.pi).floor()
15         # Finally, calculate phi using the equation in SphereFace
16         ... paper
17         phi_theta = (-1) ** (k) * cos_m_theta - 2 * k
18         # ...

```

### *Implement Loss*

```
1  def forward(self, input, target):
2      # ...
3
4      # IMPLEMENT loss
5      # Construct a tensor with index-i equals 1 otherwise 0
6      #     Initialize with all 0 in index tensor, then fill in
7      #     ... 1 in corresponding places
8      index = cos_theta * 0.0
9      index.scatter_(1, target.data.view(-1, 1), 1)
10     # Convert index to int for further use
11     index = Variable(index.to(torch.uint8))
12
13     # Notice the difference between the angular softmax and
14     # ... regular softmax
15     # Prepare the input for softmax by replace cos_function
16     # ... with phi-function
17     # at index-i (already calculated in index)
18     output = cos_theta * 1.0
19     # Here, it turns out to be crucial to use the annealing
20     # ... optimization
21     # in Appendix G in SphereFace Paper (for more
22     # ... discussion, see the assignment report)
23     # We use lambda starting from LambdaMax and gradually
24     # ... decay to LambdaMin
25     # for the decay rate we choose the one based on https
26     # ... ://github.com/clcarwin/sphereface_pytorch
27     lamb = max(self.LambdaMin, self.LambdaMax / (1+0.1*self.
28     ... it))
29     # Implement annealing
30     output[index] -= cos_theta[index] / (1. + lamb)
31     output[index] += phi_theta[index] / (1. + lamb)
32
33     # Use softmax from Pytorch for loss function
34     logpt = F.log_softmax(output)
35     # Calculate the loss functiong use eq.(7) in SphereFace
36     # ... paper
37     loss = (-logpt.gather(1, target).view(-1)).mean()
38
39     # ...
```

Here, we choose hyper-parameters such as  $k$  based on the reference of [5]. To be spe-

cific,  $k = \frac{m\theta_{y_i,i}}{\pi}$ . Another technique we use is annealing optimization in A-Softmax loss, as introduced in Appendix G in [2], where we let

$$f_{y_i} = \frac{\lambda ||\mathbf{x}_i|| \cos(\theta_{y_i}) + ||\mathbf{x}_i|| \psi(\theta_{y_i})}{1 + \lambda} \quad (3)$$

Particularly, as used in [2],  $\lambda$  starts from a large value (1500) and gradually decreased (divided by  $(1 + 0.1 \times iteration\_number)$ ) to a small value (5). The reason that this annealing strategy is essential, as described in [6], is:

”...when training data has too many subjects (such as **CASIA-WebFace dataset**), the convergence of L-Softmax will be more difficult than softmax loss.”

and this reason also applies on A-Softmax loss.

## C.2 Training Process and Evaluation Results

Now we can train the faceNet with A-Softmax loss by running *casia\_train.py*. With the help of annealing optimization, the requirement evaluation accuracy on LFW ( $> 90\%$ ) can be easily achieved just using the default arguments. Additionally, here we show 3 training process with different  $m$  and *iteration\_number* in Figure 3 and corresponding evaluation results on LFW in Table 1.

$m$ value, iterations	$m = 3$ , 28k	$m = 3$ , 40k	$m = 4$ , 28k	$m = 4$ , 40k
Accuracy	<b>0.9850</b>	0.9830	0.9757	0.9720

Table 1: Evaluation results on LFW. Evaluation Results include model trained with different  $m$  (3 vs 4) and different iteration numbers (28k vs 40k).

Here, all evaluations on LFW are have good performance ( $\geq 97.2\%$ ). One thing that is different to the experiments in [2] is that, our results show that  $m = 3$  has better evaluation accuracy than  $m = 4$ . However, this is understandable, since larger  $m$  is more difficult to train. Actually, from Figure 3 and the comparison between the evaluation accuracy of 28k and 40k training, we can see that the learning rate is too large in the last several thousand iterations, so that for  $m = 4$ , it is very possible that the training is not good enough for the used learning rate. **One thing we can try to make better performance is to use smaller learning rates and smaller decay coefficients.**

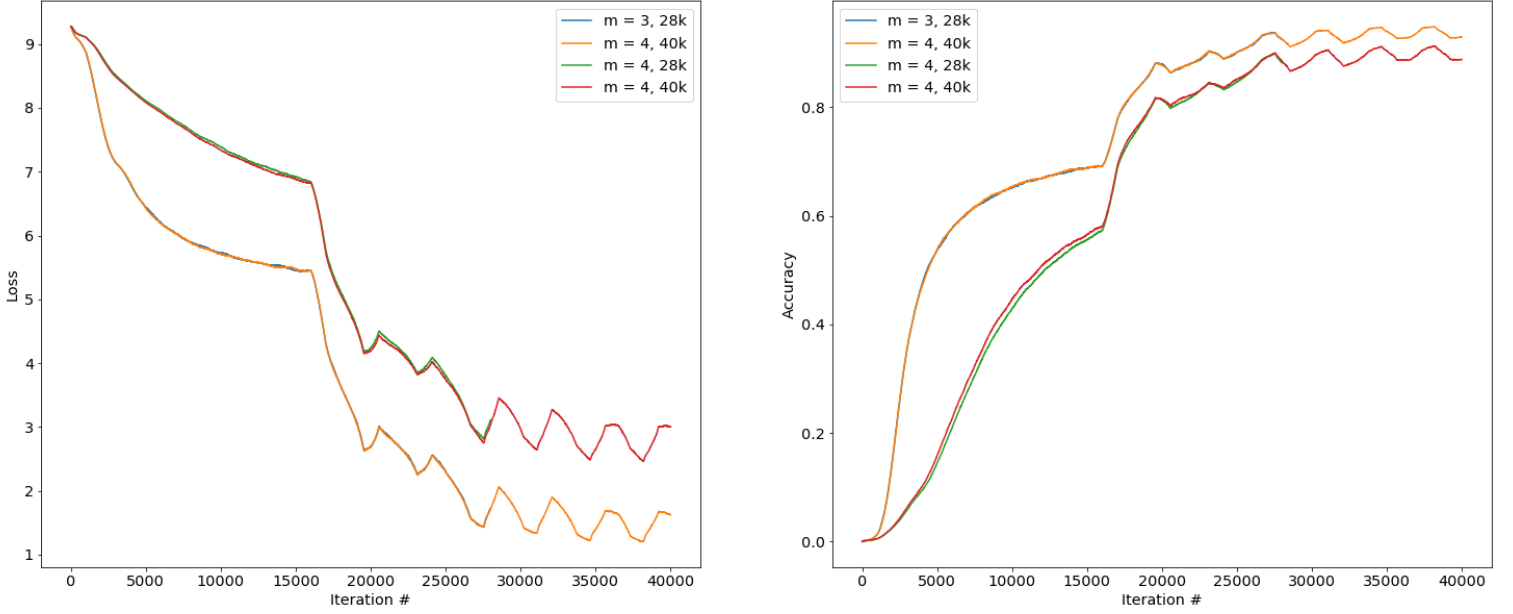


Figure 3: Training process of (a)loss and (b)accuracy of ShereFace net, including model trained with different  $m$  (3 vs 4) and different iteration numbers (28k vs 40k). **All curves are smoothed by a quadratic interloation.**

### C.3 Batch Normalization Layers

In this part, we will add batch normalization (BN) layers to the original 20-layer residual network. The BN layers will be added after each convolution layer and fully connected layer. The specific way we do this with PyTorch in code, is to both add the layer (*BatchNorm2d()* after convolutional layer and *BatchNorm1d()* after fully connected layer) in FaceNet model, and also add the explicit implementations as a layer in *forward()* step, an example for the pair described above is:

*Add BN layers in model*

```

1      self.conv1_1 = nn.Conv2d(3, 64, 3, 2, 1)  #=>B*64*56*48
2      self.bn1_1 = nn.BatchNorm2d(64)
3      self.relu1_1 = nn.PReLU(64)
4      self.conv1_2 = nn.Conv2d(64, 64, 3, 1, 1)
5      self.bn1_2 = nn.BatchNorm2d(64)
6      self.relu1_2 = nn.PReLU(64)
7      self.conv1_3 = nn.Conv2d(64, 64, 3, 1, 1)
8      self.bn1_3 = nn.BatchNorm2d(64)
9      self.relu1_3 = nn.PReLU(64)

```

### *Implement BN layers in forward() function*

```
1      x = self.relu1_1(self.bn1_1(self.conv1_1(x)))
2      x = x + self.relu1_3(self.bn1_3(self.conv1_3(self.relu1_2
      ... (self.bn1_2(self.conv1_2(x))))))
```

#### C.3.1 Training Process and Evaluation Results with Batch Normalization

With the new model that has batch normalization layers, we now retrain the network to see the influence of batch normalization. Here, we show 2 different training results with different total iteration number in Figure 4, and the corresponding evaluation results on LFW dataset in Table 2. Note that in this part, all training are based on fixed  $m$  value ( $m = 4$ ), since this part aims mainly on how the BN layers influence the model. Therefore, the best evaluation here ( $m = 4$ , w/ BN, 28k) is not better than the best evaluation in previous part ( $m = 3$ , w/o BN, 28k).

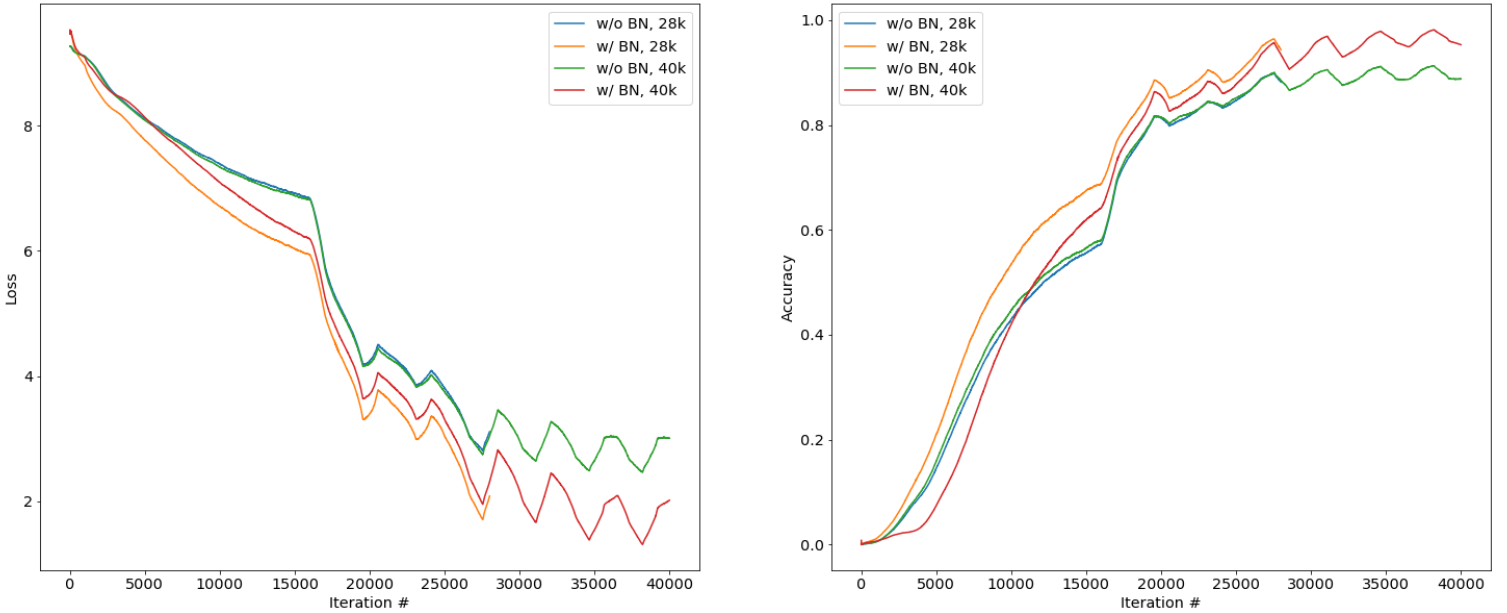


Figure 4: Training process of (a)loss and (b)accuracy of ShereFace net, including model trained with and without batch normalization (BN) layers and with different iteration numbers (28k vs 40k). **All curves are smoothed by a quadratic interpolation.**

model, iterations	w/o BN, 28k	w/ BN, 28k	w/o BN, 40k	w/ BN, 40k
Accuracy	0.9757	<b>0.9808</b>	0.9720	0.9732

Table 2: Evaluation results on LFW. Training with  $m = 4$  recommended in [2]. Evaluation Results include model trained with and without batch normalization (BN) layers and with different iteration numbers (28k vs 40k). **All curves are smoothed by a quadratic interpolation.**

### C.3.2 Comparison to Original Model

The batch normalization method is widely used in deep neural networks. The basic idea is to normalize the batch samples into a Gaussian distribution by

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{Var[x^{(k)}]}} \quad (4)$$

while also allows the network to squash the range if it wants to by

$$y^{(k)} = \gamma^{(k)}\hat{x}^{(k)} + \beta^{(k)} \quad (5)$$

where  $\gamma^{(k)} = \sqrt{Var[x^{(k)}]}$  and  $\beta^{(k)} = E[x^{(k)}]$  can be learned by the network.

From Table 2, we can see that **with BN layers, the evaluation will have better performance** both for smaller and larger number of iterations. There are several reasons that the batch normalization (BN) can help the performance of a deep neural network. First, BN layers limit the covariate shift by normalizing the activations of each layer, which results in learning on a more stable distribution of inputs at each layer, and also accelerates the training of the network. Second, BN makes the mean and variance of the activations independent of the values, which leads to the magnitude of the higher order interactions are going to be suppressed, allowing larger learning rates to be used. Third, the form of Gaussian distribution in BN can improve the gradient flow through the network. These all together makes the deep neural networks with BN layers usually outperform ones without them (Reference: [7]).

## C.4 Visualization with tSNE

One interesting result to evaluate is how the network with loss function after training can separate the different categories samples. Here to visualize the result, we will implement the tSNE on the feature vectors given by our trained network on 10 random chosen identities. We will use the tSNE codes from *sklearn* package. The script that can automatically choose random identities, feed forward the data and visualize dimension-reduced feature vector with tSNE is as following:

```

1 import torch
2 from torch.autograd import Variable
3 import faceNet, dataLoader
4 from torch.utils.data import DataLoader

```

```

5
6 import os, random
7 import numpy as np
8 import matplotlib.pyplot as plt
9 from sklearn.manifold import TSNE
10
11 def prepare_data(identities_num = 10):
12     """
13     Prepare images and corresponding landmarks of 10 random
14     ... identities
15     """
16     imageRoot = '../CASIA-WebFace/'
17     people = [f for f in os.listdir(imageRoot)]
18     identities = random.sample(people, 10)
19
20     alignmentRoot = './data/casia_landmark.txt'
21     with open(alignmentRoot) as file:
22         name = ''
23         landmarks = {}
24         for line in file:
25             new = line.split('/')[0]
26             if new == name:
27                 landmarks[name].append(line)
28             else:
29                 landmarks[new] = [line]
30                 name = new
31
32     record = []
33     identities_img_number = []
34     for i in identities:
35         identities_img_number.append(len(landmarks[i]))
36         for line in landmarks[i]:
37             record.append(line)
38     batch_size = len(record)
39
40     new_landmark = './data/casia_landmark_' + str(identities_num)
41     ... + '_identities.txt'
42     file = open(new_landmark, 'w')
43     for line in record[:-1]:
44         file.write(line)
45     file.write(record[-1][:-2])
46     file.close()
47
48     new_dic = '../CASIA-WebFace_' + str(identities_num) + '
49     ... _identities/'

```

```

46     os.system('rm -rf ' + new_dic)
47     os.system('mkdir ' + new_dic)
48
49     for i in identities:
50         os.system('cp -r ../CASIA-WebFace/' + i + ' ' + new_dic +
51             ... i)
52
53     assert(sum(identities_img_number) == batch_size)
54     return identities, identities_img_number, batch_size, new_dic
55     ... , new_landmark
56
57 identities_num = 10
58 identities, counts, batch_size, imageRoot, alignmentRoot =
59     ... prepare_data(identities_num)
60 faceDataset = dataLoader.BatchLoader(
61     imageRoot = imageRoot,
62     alignmentRoot = alignmentRoot,
63     cropSize = (96, 112)
64 )
65 faceLoader = DataLoader(faceDataset, batch_size = batch_size,
66     ... num_workers = 16, shuffle = False )
67 imBatch = Variable(torch.FloatTensor(batch_size, 3, 112, 96))
68
69 net = getattr(faceNet, 'faceNet')()
70 net.load_state_dict(torch.load("./checkpoint_bn/netFinal_8.pth")
71     ... )
72 net.eval()
73 net.feature = True
74
75 for i, dataBatch in enumerate(faceLoader):
76     image_cpu = dataBatch['img']
77     imBatch.data.resize_(image_cpu.size() )
78     imBatch.data.copy_(image_cpu )
79     output = net(imBatch)
80     break
81
82 ff = output.data.cpu().numpy()
83 visualization = TSNE(n_components=2).fit_transform(ff)
84
85 # Note that when the dataLoader.BatchLoader read in images, it
86     ... shuffles the order
87 # To visualize the tSNE result corresponding to the correct
88     ... identity, we need to

```



```

83 # restore the shuffle using the order in dataLoader.BatchLoader.
    ... perm
84 dic = {x:[] for x in range(10)}
85 perm = np.array(faceDataset.perm)
86 cum_counts = [-1] + list(np.array(counts).cumsum() - 1)
87 for i in range(identities_num):
88     dic[i] = list(np.where((perm >= cum_counts[i]+1) & (perm <=
        ... cum_counts[i+1]))[0])
89     assert(len(dic[i]) == counts[i])
90
91 figure = plt.figure(figsize=(16, 12))
92 for i in range(identities_num):
93     plt.scatter(visualization[dic[i],0], visualization[dic[i],1])
94 plt.legend(identities)
95 figure.savefig('sphere tSNE.png')

```

With the script above, we can choose then random identities for the tSNE visualization. After running several times of the experiment, we show the tSNE visualization in Figure 5, where we can clearly see that different person's feature vector will cluster separately in the tSNE reduced space.

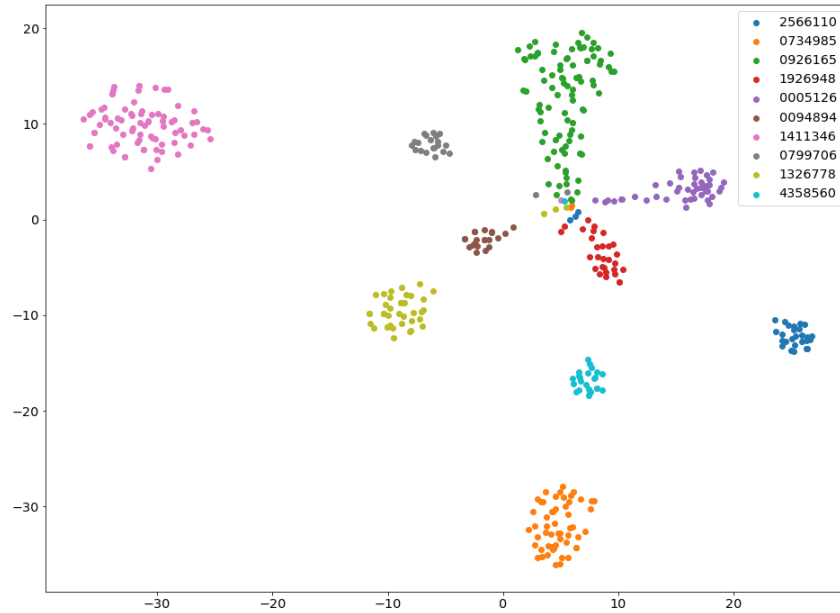


Figure 5: tSNE visualization of 10 random chosen identities, whose images are fed forward through the trained SphereFace net.

## D Training CosFace on CASIA Dataset

In this section, we will use the loss function in CosFace [8] for the similar tasks in C, including training, evaluation and visualization with tSNE.

### D.1 Cosine Loss and $\psi$ Function

Similar to the way we did in C, we first follow the cosine loss function defined by Eq.[4] in [8]:

$$L_{lmc} = \frac{1}{N} \sum_i -\log\left(\frac{e^{s\psi(\theta_{y_i,i})}}{e^{s\psi(\theta_{y_i,i})} + \sum_{j \neq y_i} e^{s\cos(\theta_{j,i})}}\right) \quad (6)$$

where

$$\psi(\theta_{y_i,i}) = \cos(\theta_{j,i}) - m \quad (7)$$

Again, in *cosface/faceNet.py*, we implement Eq.(7) under class *CustomLinear* and Eq.(6) under class *CustomLoss* respectively. The following snippets showing the the specific implementation with PyTorch:

*Implement psi\_theta*

```
1  # def forward(self, input):
2      # ...
3
4
5      # IMPLEMENT phi_theta
6      # phi = s * (cos_theta - m), in this code, s is defined
7      # ... and used in CustomLoss class
8      phi_theta = cos_theta - self.m
9
10     # ...
```

*Implement Loss*

```
1  # def forward(self, input, target):
2      # ...
3
4      # IMPLEMENT loss
5      # Construct a tensor with index-i equals 1 otherwise 0
6      # Initialize with all 0 in index tensor, then fill in
7      # ... 1 in corresponding places
8      index = cos_theta * 0.0
9      index.scatter_(1, target.data.view(-1, 1), 1)
10     # Convert index to int for further use
11     index = Variable(index.to(torch.uint8))
```

```

12     # Prepare the input for softmax by replace cos_function
    ... with phi-function
13     # at index-i (already calculated in index)
14     # at [index], output is phi_theta otherwise cos_theta,
    ... this can be simply achived use torch.where()
15     output = torch.where(index, phi_theta, cos_theta)
16     # Use softmax from Pytorch for loss function
17     logpt = F.log_softmax(self.s * output)
18     # Calculate the loss functiong use equation from CosFace
    ... paper
19     loss = (-logpt.gather(1, target).view(-1)).mean()
20
21     # ...

```

## D.2 Training Process and Evaluation Results with Cos-Loss

With the loss function defined above, we copied the exactly same network structure that we used in C, and trained the network with the Cos-Loss. Here, to compare with SphereFace fair and effectively, we train all the CosFace net here with **same learning rate and iteration numbers (28k) to the ones used in the best SphereFace net training**. The main aim this part is to play with the loss related hyperparameters ( $s$  and  $m$ ), so that in the following experiments we will fix the learning rate and iteration number, and test with different combinations of  $(s, m)$ . We choose two  $s$  values, which are the default value used in *casia\_train.py* ( $s = 30$ ) and in the *CustomLoss* class constructor( $s = 64$ ), and the choices of  $m$  are based on [9]. The training process is shown in Figure 6, and the the corresponding evaluation results are shown in Table 3.

$s, m$	$s = 30, m = 0.35$	$s = 30, m = 0.40$	$s = 64, m = 0.35$	$s = 64, m = 0.40$
Accuracy	<b>0.9902</b>	0.9882	0.9868	0.9885

Table 3: Evaluation results of FaceNet with Cos-Loss. All training process has same learning rate and iteration number to the best SphereFace net we have in Sec. C

## D.3 Compared to SphereFace

Compare the best results we have from the CosFace and SphereFace nets, we can see that the CosFace outperforms the SphereFace by a margin about 1%. Since we makes the training procedure to be same in both cases, the better performance is due to the difference between A-Softmax and Cosine Loss function, as discussed in [8, 2].

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

$$C_1 : \cos(m\theta_1) \geq \cos(\theta_2)$$

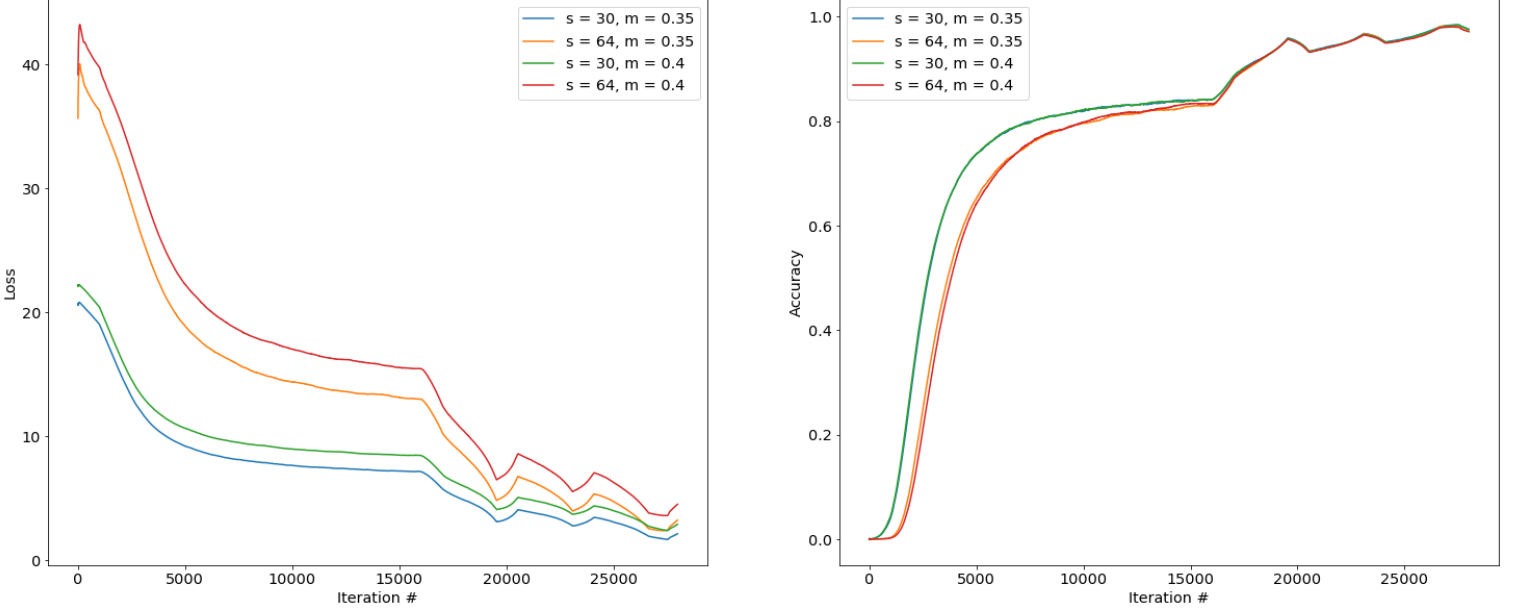


Figure 6: Training process of (a)loss and (b)accuracy of CosFace net, including model trained with different combinations of  $s$  and  $m$ . **All curves are smoothed by a quadratic interpolation.**

$$C_1 : \cos(m\theta_2) \geq \cos(\theta_1)$$

thus, the decision margin that classifies the two classes are  $\theta_1 \leq \frac{\theta_2}{m}$  for  $C_1$  and  $\theta_2 \leq \frac{\theta_1}{m}$  for  $C_2$ .

On the other hand, for Cosine loss, the decision margin is defined in the cosine space rather than the angle space like A-softmax, thus, the boundary is given by:

$$C_1 : \cos(\theta_1) \geq \cos(\theta_2) + m$$

$$C_1 : \cos(\theta_2) \geq \cos(\theta_1) + m$$

Therefore,  $\cos(\theta_1)$  is maximized while  $\cos(\theta_2)$  being minimized for  $C_1$  and similarly for  $C_2$  to perform large-margin classification, and as a result, a clear margin ( $\sqrt{2}m$ ) in the produced distribution of the consine of angle. An illustrative figure to compare A-softmax and Cosine loss is shown in Figure 7. Consequently, given the same network structure, using Cosine loss function can be more robust than A-softmax and yield a better evaluation result on LFW dataset. Actually, this is consistent to those results achieved in [8].

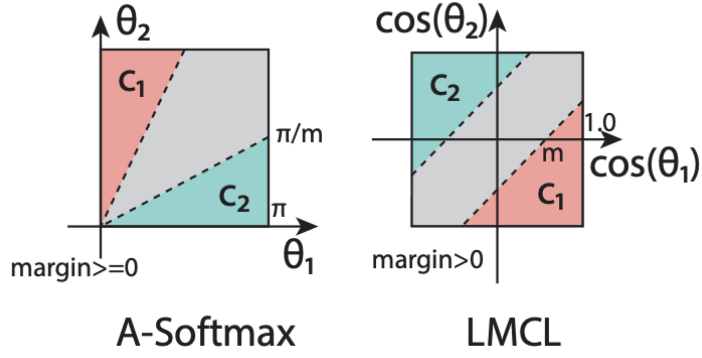


Figure 7: The comparison of decision margins between A-Softmax and Large Margin Cosine Loss, for the binary-classes scenarios. Dashed line represents decision boundary, and gray areas are decision margins. Figure borrowed from [8]

#### D.4 Visualization with tSNE

In this last part, we will use the tSNE again to visualize the feature vectors of the 10 same identities' images. This time, those images will be fed forward to the best CosFace net we trained. The results are shown in Figure 8. We can see that the feature vector clustered more obviously than the previous SphereFace net results.

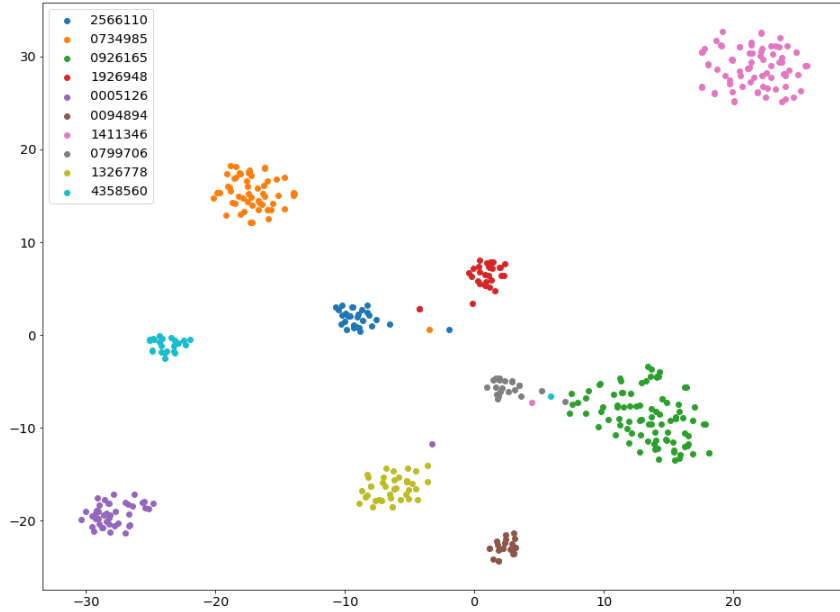


Figure 8: tSNE visualization of 10 same identities chosen in the previous section, whose images are fed forward through the trained CosFace net.

## References

- [1] 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. Technical Report 07-49, University of Massachusetts, Amherst, October 2007.
- [2] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphreface: Deep hypersphere embedding for face recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 212–220, 2017.
- [3] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. IEEE Signal Processing Letters, 23(10):1499–1503, 2016.
- [4] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. arXiv preprint arXiv:1411.7923, 2014.
- [5] [https://github.com/clcarwin/sphreface\\_pytorch](https://github.com/clcarwin/sphreface_pytorch)
- [6] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. Large-Margin Softmax Loss for Convolutional Neural Networks. In ICML, vol. 2, no. 3, p. 7. 2016.

- [7] Ioffe, Sergey, and Christian Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift." arXiv preprint arXiv:1502.03167 (2015).
- [8] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 5265–5274, 2018.
- [9] [https://github.com/MuggleWang/CosFace\\_pytorch](https://github.com/MuggleWang/CosFace_pytorch)