



# CS F425: Deep Learning

# 13

## Loss Function For Neural NW



**Dr. Kamlesh Tiwari**  
Assistant Professor, Department of CSIS,  
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

Feb 17, 2023 **ON-CAMPUS** Campus @ BITS-Pilani [Jan-May 2023]

<http://ktiwari.in/dl>

### Lets look closer

- Consider two class classification and a single example
- Cross Entropy** loss is:

$$\begin{aligned}
 L &= \frac{1}{m} \sum_{i=1}^C -y_i \log(f(s_i)) \\
 &= \sum_{i=1}^2 -y_i \log(f(s_i)) = -y_1 \log(f(s_1)) - y_2 \log(f(s_2)) \\
 &= -y_1 \log(f(s_1)) - (1 - y_1) \log(1 - f(s_1))
 \end{aligned}$$

- Output vector is *one-hot* for the multi-class classification, so most of the  $y_i$  are zero except a single target positive one, leading to

$$L = -\log\left(\frac{e^{s_p}}{\sum_{j=1}^C e^{s_j}}\right)$$

here  $s_p$  is output with respect to the positive class of the input.

### Centre Loss

- Calculate center<sup>1</sup> for each class, let  $c_{y_i}$  is centre of deep features belonging to  $y_i^{th}$  class
- Then move inter-correlated features closer to the centre of the class

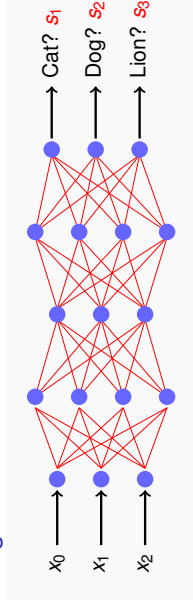
$$L_c = \frac{1}{2} \sum_m \|x_i - c_{y_i}\|_2^2$$

### Issues

- Determining center is expensive
- Centre is not accurate while using batch
- Euclidean distance is not the best measure for the feature similarity

<sup>1</sup>average of features of a class

### Training Loss for Multi-Class classification



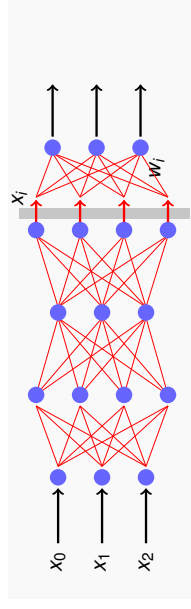
- Softmax** can help make summation one (for being probability)

$$f(s_i) = \frac{e^{s_i}}{\sum_{j=1}^C e^{s_j}}$$

- Loss** can tell how good these values are  
**Cross Entropy** is the right choice:

$$L = \frac{1}{m} \sum_{i=1}^C -y_i \log(f(s_i))$$

### The Cross Entropy loss



$$L = -\log\left(\frac{e^{s_p}}{\sum_{j=1}^C e^{s_j}}\right) = -\frac{1}{m} \sum_{i=1}^m \log\left(\frac{e^{W_i^T x_i + b_i}}{\sum_{j=1}^C e^{W_j^T x_i + b_j}}\right) \quad (1)$$

- where  $x_i$  is deep feature of  $i$ th sample belonging to class  $y_i$

Softmax loss is separable but not discriminative enough

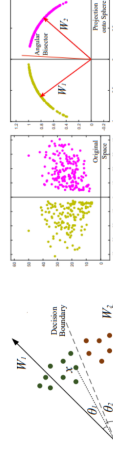
### Cosine Formulation of the Softmax Loss

- $W_j^T x_i + b_j$  can be taken as  $W_j^T x_i$ , which is equal to

$$\|W_j\| \cdot \|x_i\| \cos \theta_{j,i}$$

- If we normalize the weights, making  $\|W_j\| = 1$   
Modified loss function is represented as

$$L = -\frac{1}{m} \sum_i \log\left(\frac{e^{\|x_i\| \cos \theta_{y_i,i}}}{\sum_j e^{\|x_i\| \cos \theta_{j,i}}}\right)$$



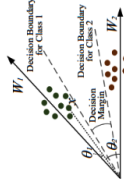
Loss is angularly distributed

## ShpereFace: angular softmax

Incorporate multiplicative angular margin  $m \geq 1$

$$L = -\frac{1}{m} \sum_i \log \left( \frac{e^{||x_i||\psi(\theta_{y_i,i})}}{e^{||x_i||\psi(\theta_{y_i,i})} + \sum_{j \neq y_i} e^{||x_i||\cos(\theta_{j,i})}} \right)$$

- $\psi(\theta_{y_i,i}) = (-1)^k \cos(m\theta_{y_i,i}) - 2k$
- $\theta_{y_i,i} \in \left[ \frac{k\pi}{m}, \frac{(k+1)\pi}{m} \right]$
- $k \in [0, m]$



Margin enforce compression for intra class feature distribution and expands inter-class margin.

**Issues:** different margin for different classes, difficult to train.

## CosFace: Large Margin Cosine Loss <sup>2</sup>

- 1 Introduces margin  $m \geq 0$  in cosine difference between classes

$$L = -\frac{1}{m} \sum_i \log \frac{e^{s(\cos(\theta_{y_i,i})-m)}}{e^{s(\cos(\theta_{y_i,i})-m)} + \sum_{j \neq y_i} e^{s \cdot \cos(\theta_{j,i})}}$$

- 2 Defines a decision margin in cosine space rather than the angle space
- 3 **Issue:** nonlinear angular margin

<sup>2</sup>Wang, Hao, et al. "Cosface: Large margin cosine loss for deep face recognition." Proceedings of the IEEE conference on computer vision and pattern recognition. 2018.

Thank You!

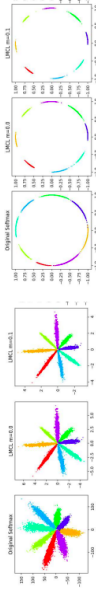
Thank you very much for your attention! <sup>4</sup>

## Hypersphere

- 1 With feature similarity as cosine,  $||x_j||$  does not contribute to score
- 2 Let us fix it  $||x_j|| = s$
- 3 And use loss as

$$L = -\frac{1}{m} \sum_i \log \left( \frac{e^{s \cdot \cos(\theta_{y_i,i})}}{\sum_j e^{s \cdot \cos(\theta_{j,i})}} \right)$$

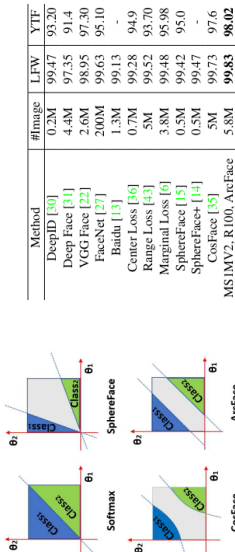
- 4 Feature vectors are distributed on hypersphere of radius  $s$ .



**ArcFace:** <sup>3</sup>  
Introduces margin  $m \geq 1$  to the classification boundary

$$L = -\frac{1}{m} \sum_i \log \frac{e^{s \cdot \cos(\theta_{y_i,i}+m)}}{e^{s \cdot \cos(\theta_{y_i,i}+m)} + \sum_{j \neq y_i} e^{s \cdot \cos(\theta_{j,i})}}$$

$||w_j|| = 1$ , and  $x_i$  is  $L_2$  normalized and scaled to  $s$



<sup>3</sup>Deng, Jiarui, et al. "Arcface: Additive angular margin loss for deep face recognition." Proceedings of the IEEECVF Conference on Computer Vision and Pattern Recognition. 2019.

<sup>4</sup>Adopted from:

<https://towardsdatascience.com/additive-margin-softmax-loss-am-softmax-912e11ce1c6b>