

# CNN Classification and Verification

Hao Chen

# Content

- Basics
- Face Classification and Verification
- How to Learn Better Features?
- Augmentation & Regularization

# Basics

# Classification

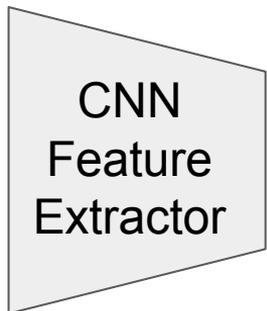
Images

Features

Logits

Probabilities

$$I \in \mathbb{R}^{H \times W \times 3}$$



$$\mathbf{f} \in \mathbb{R}^D$$

$$W \in \mathbb{R}^{C \times D}$$



$$\mathbf{z} \in \mathbb{R}^C$$

$$\begin{bmatrix} W_0 \\ W_1 \\ \dots \\ W_C \end{bmatrix}$$

Softmax

$$\mathbf{p} \in \mathbb{R}^C$$

[0.03 0.09 0.23 0.62 0.03]

Prediction

$$\hat{y} = \arg \max_{i \in \{C\}} p_i$$

# How to Train Classification?

Cross-Entropy  $H(p, q) = - \sum_x p(x) \log q(x)$

# How to Train Classification?

Cross-Entropy  $H(p, q) = - \sum_x p(x) \log q(x)$

In classification,  $p(x)$  is one-hot:  $[0, 0, 0, 1, 0]$

$q(x)$  is predicted probability from softmax  $\left[ \frac{e^{W_0^T \mathbf{f}}}{\sum_{j=1}^C e^{W_j^T \mathbf{f}}}, \frac{e^{W_1^T \mathbf{f}}}{\sum_{j=1}^C e^{W_j^T \mathbf{f}}}, \dots, \frac{e^{W_C^T \mathbf{f}}}{\sum_{j=1}^C e^{W_j^T \mathbf{f}}} \right]$

# How to Train Classification?

Cross-Entropy  $H(p, q) = - \sum_x p(x) \log q(x)$

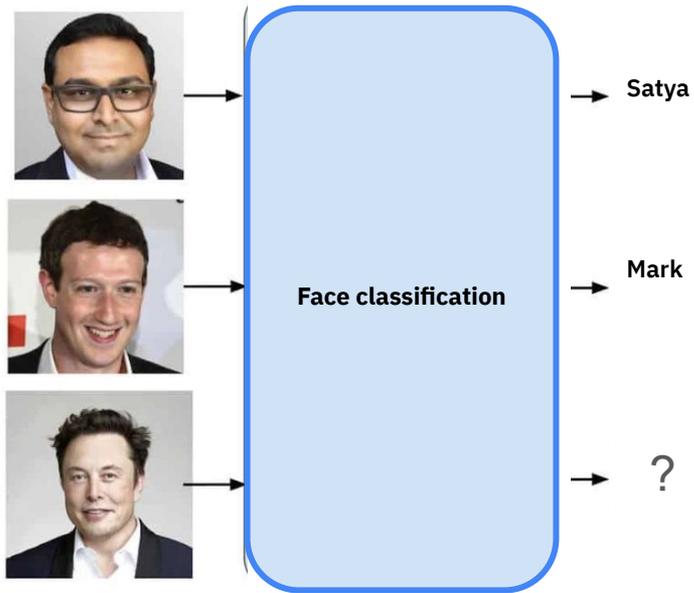
In classification,  $p(x)$  is one-hot:  $[0, 0, 0, 1, 0]$

$q(x)$  is predicted probability from softmax  $\left[ \frac{e^{W_0^T \mathbf{f}}}{\sum_{j=1}^C e^{W_j^T \mathbf{f}}}, \frac{e^{W_1^T \mathbf{f}}}{\sum_{j=1}^C e^{W_j^T \mathbf{f}}}, \dots, \frac{e^{W_C^T \mathbf{f}}}{\sum_{j=1}^C e^{W_j^T \mathbf{f}}} \right]$

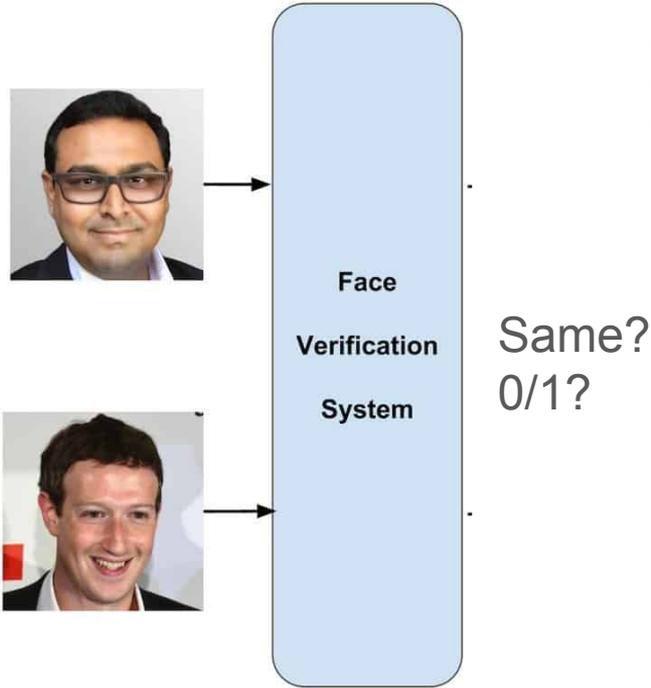
Cross-Entropy Loss:  $-\frac{1}{N} \sum_{i=1}^N \log \frac{e^{W_{y_i}^T \mathbf{f}_i}}{\sum_{j=1}^C e^{W_j^T \mathbf{f}_i}}$

# Face Classification and Verification

# Face Classification



# Face Verification



This is usually achieved by **comparing features of two faces**

# Distance vs Similarity between Vectors

How to compare two face features?

Distance:  $\|\mathbf{f}_1 - \mathbf{f}_2\|^2 = \|\mathbf{f}_1\|^2 + \|\mathbf{f}_2\|^2 - 2\mathbf{f}_1\mathbf{f}_2^T$

Cosine Similarity:  $\cos(\mathbf{f}_1, \mathbf{f}_2) = \frac{\mathbf{f}_1\mathbf{f}_2}{\|\mathbf{f}_1\|\|\mathbf{f}_2\|}$

# Why Classification Network is not Good for Verification?

Features learned by the classifier with Cross-Entropy only is **not discriminative** enough

# A Closer Look at Cross-Entropy

$$\begin{aligned} & -\log \frac{\exp(z^y)}{\sum_{j=1}^C \exp(z^k)} \\ & = \log \frac{\sum_{j=1}^C \exp(z^k)}{\exp(z^y)} \end{aligned}$$

# A Closer Look at Cross-Entropy

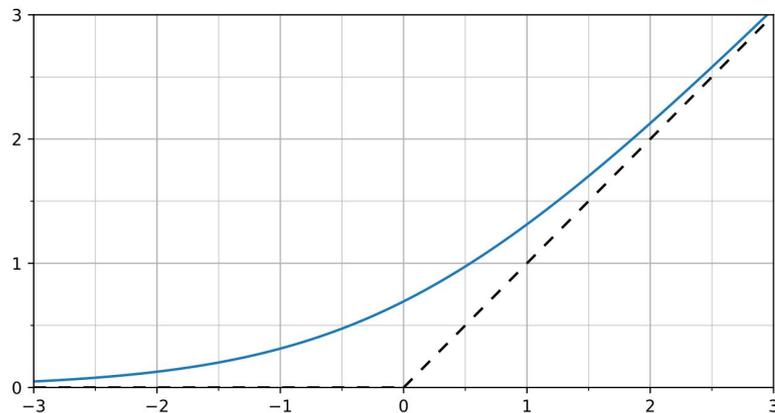
$$\begin{aligned} & -\log \frac{\exp(z^y)}{\sum_{j=1}^C \exp(z^j)} \\ &= \log \frac{\sum_{j=1}^C \exp(z^j)}{\exp(z^y)} \\ &= \log \left( 1 + \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \end{aligned}$$

# A Closer Look at Cross-Entropy

$$\begin{aligned} & -\log \frac{\exp(z^y)}{\sum_{j=1}^C \exp(z^j)} \\ &= \log \frac{\sum_{j=1}^C \exp(z^j)}{\exp(z^y)} \\ &= \log \left( 1 + \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \\ &= \log \left( 1 + \exp \left( \log \left( \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \right) \right) \end{aligned}$$

# A Closer Look at Cross-Entropy

$$\begin{aligned} & -\log \frac{\exp(z^y)}{\sum_{j=1}^C \exp(z^j)} \\ &= \log \frac{\sum_{j=1}^C \exp(z^j)}{\exp(z^y)} \\ &= \log \left( 1 + \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \\ &= \log \left( 1 + \exp \left( \log \left( \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \right) \right) \\ &\geq \max \left( \log \left( \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right), 0 \right) \end{aligned} \quad \text{Softplus } \log(1 + \exp(x)) \geq \max(x, 0)$$



# A Closer Look at Cross-Entropy

$$\begin{aligned} & -\log \frac{\exp(z^y)}{\sum_{j=1}^C \exp(z^j)} \\ &= \log \frac{\sum_{j=1}^C \exp(z^j)}{\exp(z^y)} \\ &= \log \left( 1 + \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \\ &= \log \left( 1 + \exp \left( \log \left( \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \right) \right) \\ &\geq \max \left( \log \left( \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right), 0 \right) \quad \text{Softplus } \log(1 + \exp(x)) \geq \max(x, 0) \\ &\geq \log \sum_{j=1, j \neq y}^C \exp(z^j) - z^y \end{aligned}$$

# A Closer Look at Cross-Entropy

$$\begin{aligned} & -\log \frac{\exp(z^y)}{\sum_{j=1}^C \exp(z^j)} \\ &= \log \frac{\sum_{j=1}^C \exp(z^j)}{\exp(z^y)} \\ &= \log \left( 1 + \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \\ &= \log \left( 1 + \exp \left( \log \left( \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \right) \right) \\ &\geq \max \left( \log \left( \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right), 0 \right) \quad \text{Softplus } \log(1 + \exp(x)) \geq \max(x, 0) \\ &\geq \log \sum_{j=1, j \neq y}^C \exp(z^j) - z^y \\ &\approx \max_{j \in [C], j \neq y} (z^j) - z^y \quad \text{LogSumExp } \log \sum \exp(x) \approx \max(x) \end{aligned}$$

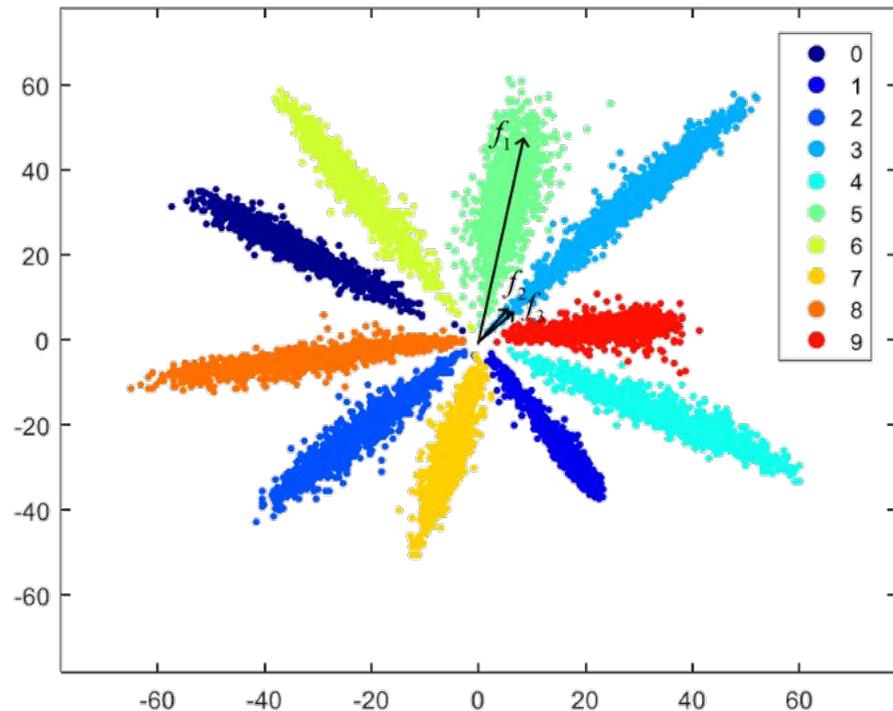
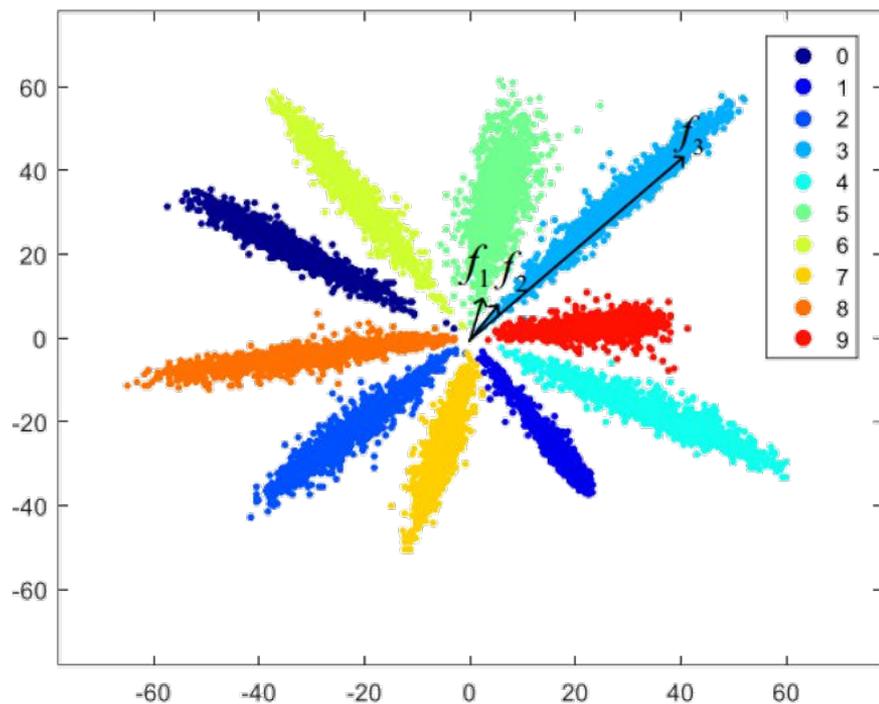
# A Closer Look at Cross-Entropy

$$\begin{aligned} & -\log \frac{\exp(z^y)}{\sum_{j=1}^C \exp(z^j)} \\ &= \log \frac{\sum_{j=1}^C \exp(z^j)}{\exp(z^y)} \\ &= \log \left( 1 + \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \\ &= \log \left( 1 + \exp \left( \log \left( \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right) \right) \right) \\ &\geq \max \left( \log \left( \sum_{j=1, j \neq y}^C \exp(z^j - z^y) \right), 0 \right) \quad \text{Softplus } \log(1 + \exp(x)) \geq \max(x, 0) \\ &\geq \log \sum_{j=1, j \neq y}^C \exp(z^j) - z^y \\ &\approx \max_{j \in [C], j \neq y} (z^j) - z^y \quad \text{LogSumExp } \log \sum \exp(x) \approx \max(x) \end{aligned}$$

Cross-Entropy optimizes

**the target logit score to be larger than the maximum of remaining logit scores**

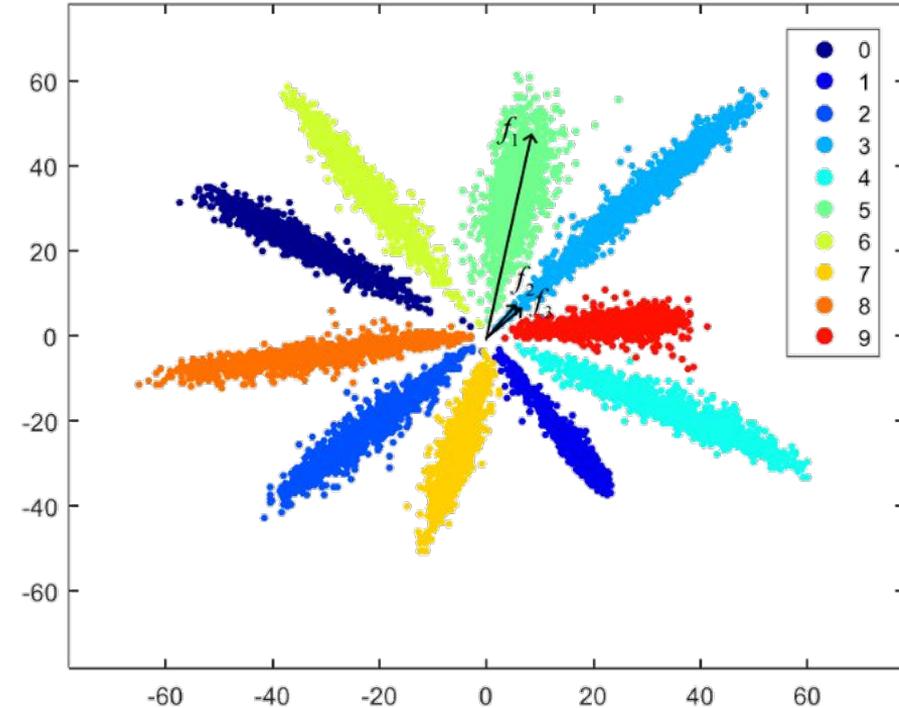
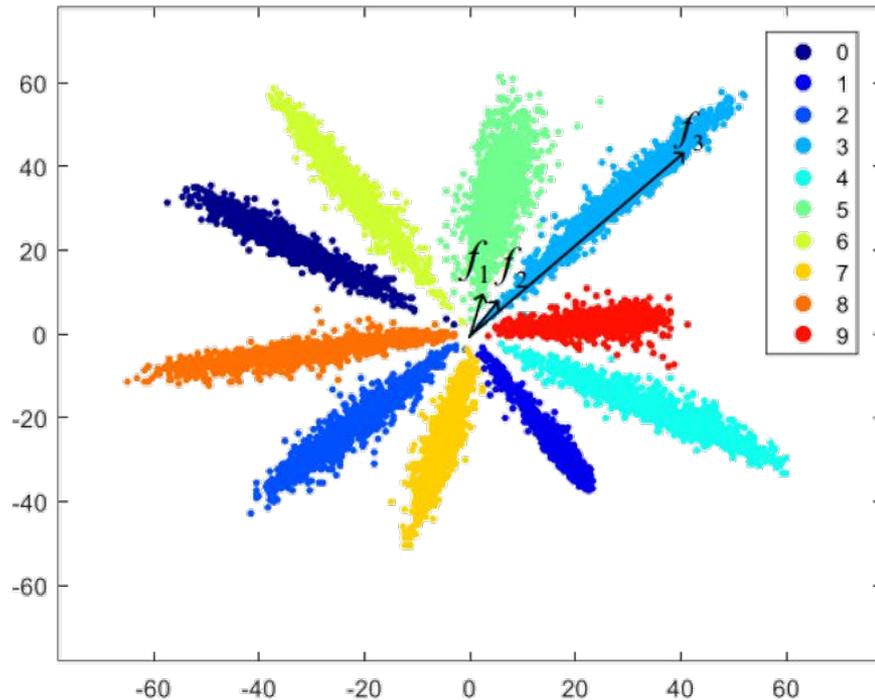
# Feature Space



Cross-Entropy leads to radius feature representations:

- It strengthens the feature length/magnitude of correctly classified samples

# Problem of Radius Feature Space



Left: Distance of features from different classes are smaller than that of same class

Right: Cosine similarity of features from different classes are larger than that of the same class

How to Learn Better Features?

# What are Good Features for Verification?

- Features of the same classes have larger similarity than other classes
  - Ideally with some margin

# What are Good Features for Verification?

- Features of the same classes have larger similarity than other classes
  - Ideally with some margin
- Optimization objectives:
  - Maximize the intra-class (same class) similarity
  - Minimize the inter-class (different class) similarity
  - Encourage a margin between intra-class and inter-class similarity

# What are Good Features for Verification?

- Features of the same classes have larger similarity than other classes
  - Ideally with some margin
- Optimization objectives:
  - Maximize the intra-class (same class) similarity
  - Minimize the inter-class (different class) similarity
  - Encourage a margin between intra-class and inter-class similarity
- Loss objective:  $z^n - z^p + m$
- But...how to define  $z^n$  and  $z^p$

# Two Paradigms

- Margin-based Softmax Losses
  - Spheraface
  - Cosface
  - Arcface
  - ...
  
- Metric/Pair-based Losses
  - Triplet Loss
  - N-Pair Loss
  - ...

# Margin-based Softmax

In Cross-Entropy, Weight vector  $W_j$  of Weight matrix  $W$  can be viewed as anchor for each class

# Margin-based Softmax

In Cross-Entropy, Weight vector  $W_j$  of Weight matrix  $W$  can be viewed as anchor for each class

Intra-class logit score:  $z^y = W_y^T \mathbf{f}$

Inter-class logit score:  $z^j = W_j^T \mathbf{f}, j \neq y$

# Margin-based Softmax

In Cross-Entropy, Weight vector  $W_j$  of Weight matrix  $W$  can be viewed as anchor for each class

Intra-class logit score:  $z^y = W_y^T \mathbf{f}$

Inter-class logit score:  $z^j = W_j^T \mathbf{f}, j \neq y$

Cross-Entropy with normal Softmax is optimizing  $\max_{j \in [C], j \neq y} (z^j) - z^y$

What is missing?

# Margin-based Softmax

In Cross-Entropy, Weight vector  $W_j$  of Weight matrix  $W$  can be viewed as anchor for each class

Intra-class logit score:  $z^y = W_y^T \mathbf{f}$

Inter-class logit score:  $z^j = W_j^T \mathbf{f}, j \neq y$

Cross-Entropy with normal Softmax is optimizing  $\max_{j \in [C], j \neq y} (z^j) - z^y$

What is missing? **Margin!**

# How to Introduce Margin in Softmax

- Normalize the feature and weight vectors -> dot product becomes cosine angle

$$\frac{e^{W_y^\top f}}{\sum_{j=1}^C e^{W_j^\top f}} = \frac{e^{\cos \theta_y}}{\sum_{j=1}^C e^{\cos \theta_j}}$$

# How to Introduce Margin in Softmax

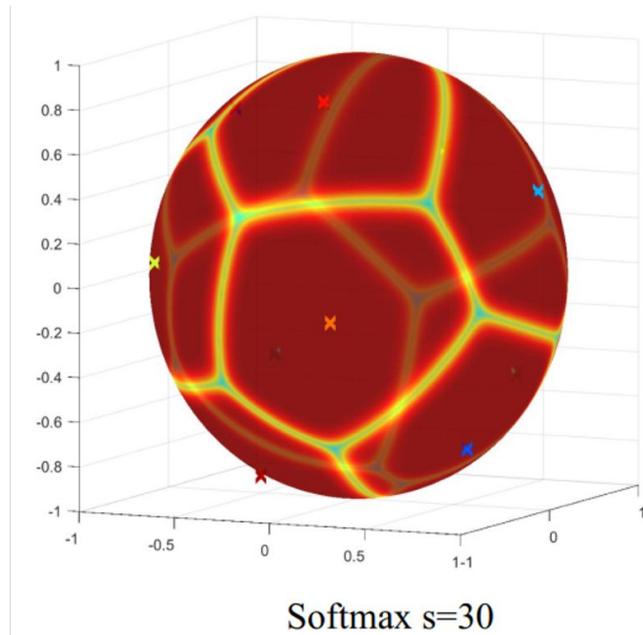
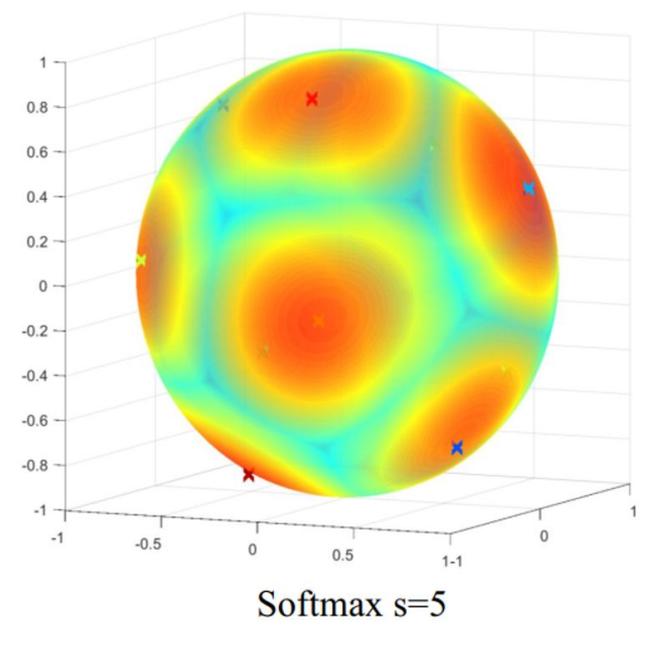
- Normalize the feature and weight vectors -> dot product becomes cosine angle

$$\frac{e^{W_y^\top f}}{\sum_{j=1}^C e^{W_j^\top f}} = \frac{e^{\cos \theta_y}}{\sum_{j=1}^C e^{\cos \theta_j}}$$

- Scale the feature to overcome the optimization problem with logit score  $[-1, 1]$

$$\frac{e^{s \cos \theta_y}}{\sum_{j=1}^C e^{s \cos \theta_j}}$$

# Effect of Scale $s$



# How to Introduce Margin in Softmax

- Normalize the feature and weight vectors -> dot product becomes cosine angle

$$\frac{e^{W_y^\top f}}{\sum_{j=1}^C e^{W_j^\top f}} = \frac{e^{\cos \theta_y}}{\sum_{j=1}^C e^{\cos \theta_j}}$$

- Scale the feature to overcome the optimization problem with logit score [-1, 1]

$$\frac{e^{s \cos \theta_y}}{\sum_{j=1}^C e^{s \cos \theta_j}}$$

- Introduce margin...at where?

# CosFace Softmax Loss

Additive Margin

$$-\log \frac{e^{s(\cos \theta_y - m)}}{e^{s(\cos \theta_y - m)} + \sum_{j \neq y}^C e^{s \cos \theta_j}}$$

# SphereFace Softmax Loss

Angular Margin

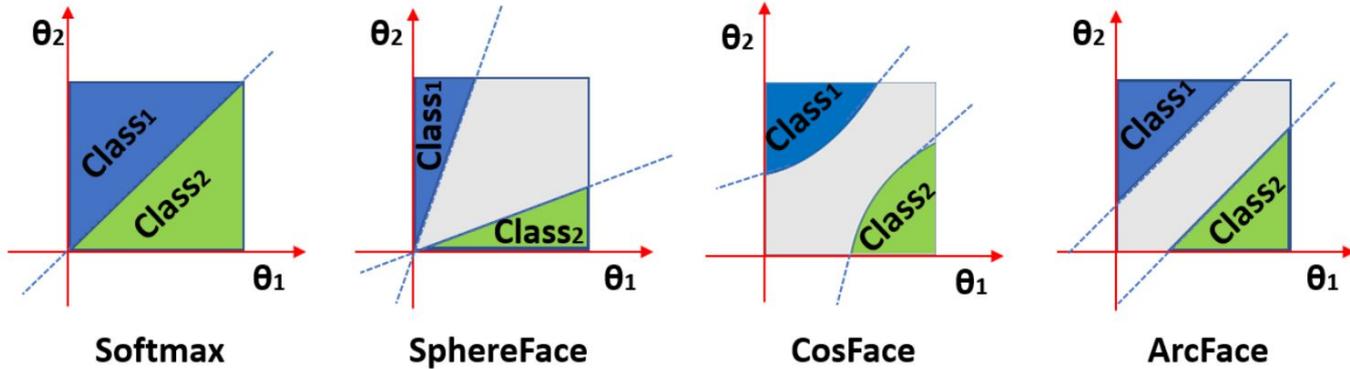
$$-\log \frac{e^{s \cos m\theta_y}}{e^{s \cos m\theta_y} + \sum_{j \neq y}^C e^{s \cos \theta_j}}$$

# ArcFace Softmax Loss

Additive angular Margin

$$-\log \frac{e^{s \cos(\theta_y - m)}}{e^{s \cos(\theta_y - m)} + \sum_{j \neq y}^C e^{s \cos \theta_j}}$$

# Margin Visualization



Different types of margin can be combined -> CombinedMarginFace

CombinedMargin Face

# Metric/Pair-based Loss

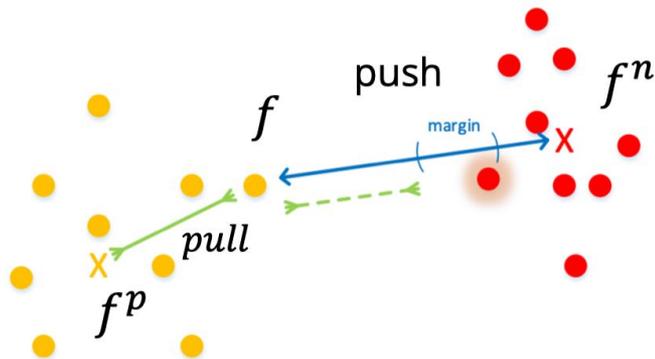
Based on pairs...with **Data Sampler**

Intra-class logit score:  $z^p = \mathbf{f}_p^T \mathbf{f}$

Inter-class logit score:  $z^n = \mathbf{f}_n^T \mathbf{f}$

# Triplet Loss

$$\mathcal{L}_{\text{Triplet}} = \max(\mathbf{f}_n^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f} + m, 0)$$



←: maximize      →: minimize

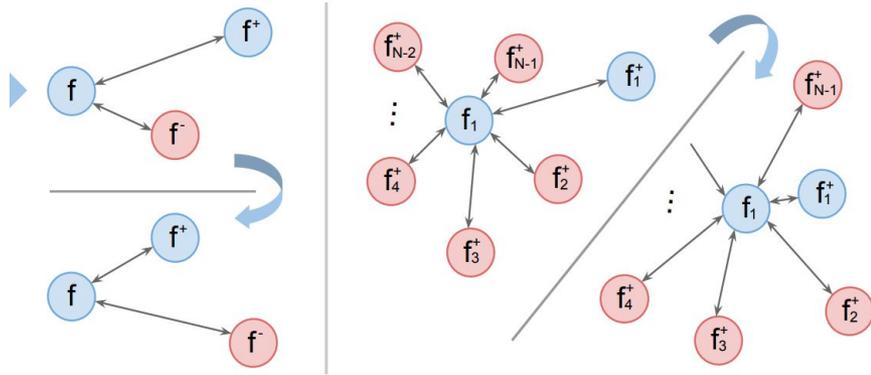
# How to Construct Positive and Negative Pairs?

- Simple way: construct from data batch
  - Might not have samples from the same classes
- Use Data Sampler to construct data batch with both samples of positive classes and negative classes
- **Hard mining**: find the most difficult pairs

# N-Pair Loss/Contrastive Loss

Why Triplet Loss is not good enough?

Only one positive pair and one negative pair



$$\log \left( 1 + \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \right)$$

## More on N-Pair/Contrastive Loss

$$\begin{aligned} & \log \left( 1 + \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \right) \\ &= \log \left( 1 + \exp \log \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \right) \end{aligned}$$

## More on N-Pair/Contrastive Loss

$$\begin{aligned} & \log \left( 1 + \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \right) \\ &= \log \left( 1 + \exp \log \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \right) \\ &\geq \max \left( \log \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}), 0 \right) \end{aligned}$$

## More on N-Pair/Contrastive Loss

$$\begin{aligned} & \log \left( 1 + \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \right) \\ &= \log \left( 1 + \exp \log \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \right) \\ &\geq \max \left( \log \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}), 0 \right) \\ &\geq \log \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \end{aligned}$$

## More on N-Pair/Contrastive Loss

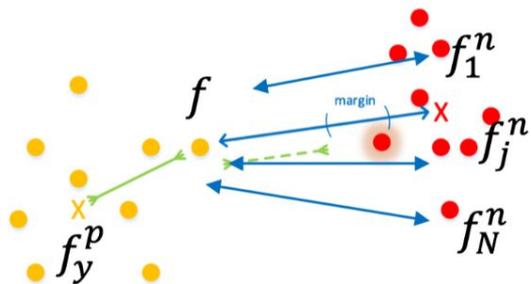
$$\begin{aligned} & \log \left( 1 + \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \right) \\ &= \log \left( 1 + \exp \log \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \right) \\ &\geq \max(\log \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}), 0) \\ &\geq \log \sum_{i=1}^{N-1} \exp(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \\ &\geq \max(\mathbf{f}_{n_i}^T \mathbf{f} - \mathbf{f}_p^T \mathbf{f}) \end{aligned}$$

- Optimization goal the same as Cross-Entropy
- But negative pairs are different
- Margin can also be added

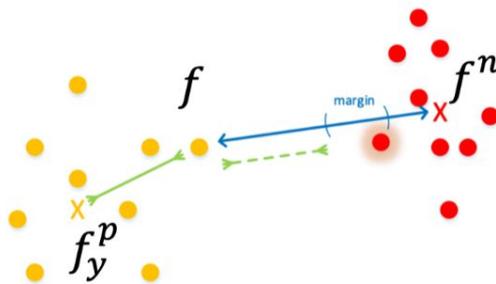
# A Unified View of Both Paradigm

Can we introduce multiple positive pairs and multiple negative pairs?

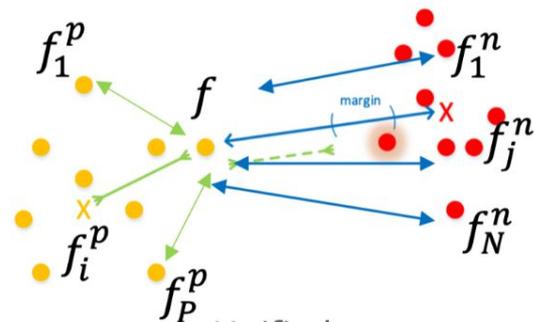
$$\log \left( 1 + \sum_i^P \sum_j^N e^{\gamma(z_{n_j} - z_{p_i} + m)} \right)$$



Am-Softmax  
# 1 positive pair  
# (C-1) negative pairs



Triplet  
# 1 positive pair  
# 1 negative pair



Unified  
# P positive pairs  
# N negative pairs

# Unified Formulation

- This is a unified formulation

$$\log \left( 1 + \sum_i^P \sum_j^N e^{\gamma(z_{n_j} - z_{p_i} + m)} \right)$$

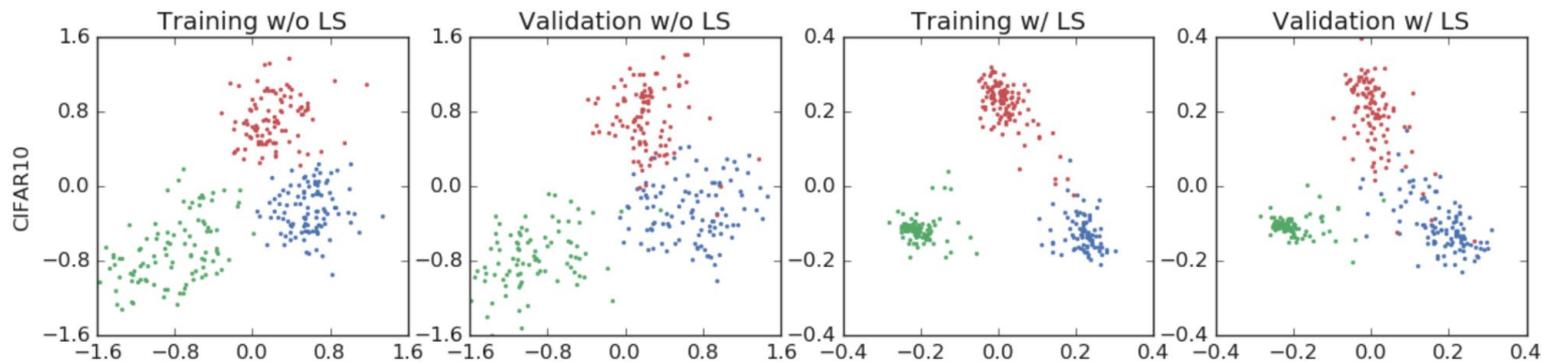
- You can derive softmax, arcface, triplet, n-pair losses from this formulation
  - By setting the number of positives/negatives and how to construct them
- **Circle loss** can also be derived from this function
- **Supervised Contrastive loss** is another loss function with multiple p/ns

# Augmentation & Regularization

# Label Smoothing

One-hot target: [0, 0, 0, 1, 0]

Smoothed ( $\text{ls}=0.1$ ) target: [0.025, 0.025, 0.025, 0.9, 0.025]



Label smoothing learns more separable features with smaller feature norm

Similar to the effect of  $s$  in Softmax

# Mixup/Cutout/Cutmix

Mixup [48]



Cutout [3]



CutMix

