# Information Amount-Guide Angular Margin Loss (IGAM Loss)
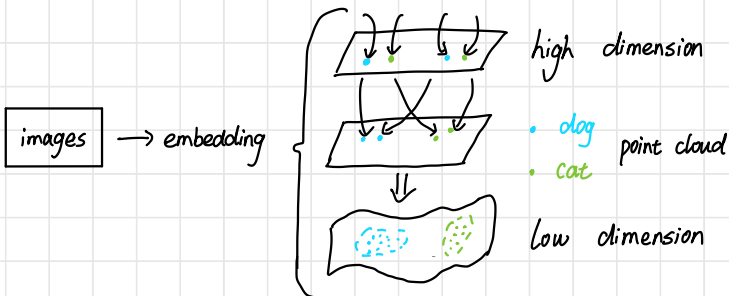
**1: Category information amount**

① Definition

Since Manifold Distribution Hypothesis:



images → embedding

high dimension

• dog
• cat
point cloud

low dimension

$\Sigma(X_i)$ determinant → Manifold Volume

$I_i = Vol(X_i) = \frac{1}{2}\log_2 \det(\Sigma(X_i)+I)$

numerical stability

positive-definite

Manifold Volume → Category information amount

② Measure

Category $i$ → information amount $I_i$ → $I_i = Vol(X_i)$

$X_i = [X_1, X_2, \cdots, X_m]$, $m$: instances number

embeddings set $X_i = [X_1, X_2, \cdots, X_j, \cdots X_m] \in R^{p\times m}$

$X_j \in R^p$ → embedding dimension

High dimension embedding space distribution characteristics → Covariance Matrix $\Sigma(X_i)$

$\Sigma(X_i) = \frac{1}{m}\sum_{j=1}^{m}(X_j-\bar{X})(X_j-\bar{X})^T$, $\bar{X} = \frac{1}{m}\sum_{j=1}^{m}X_j$

for Covariance Matrix estimation accuracy ↑

Employ Ledoit-péché nonlinear shrinkage

$\Sigma(X_i) = V diag(\lambda_1, \lambda_2, \cdots, \lambda_p) V^T$

$V$: eigenvectors Matrix $\lambda_i = max(\lambda_i, \lambda_-)$

$\lambda_-$: nonlinearly transformed minimum eigenvalue

$\lambda_- = (1-\sqrt{\frac{p}{m}})^2$
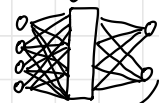
**2: IGAM Loss**

① Etiology (motivation)

classification:

images

↓

feature extractor

↓

Classify
(fcl)

$logitis_i = w_i^T x$

$W = [w_1, \cdots, w_c]$

feature vector: $x$
label: $i$

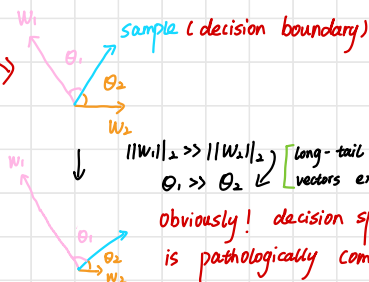Cross-entropy loss: $L = -\log(\frac{e^{w_i^T x}}{\sum_{j=1}^{c}e^{w_j^T x}})$ (softmax)

For example: Binary classification

Core: for sample locate on decision boundary

we assume $w_1^T x = w_2^T x$

$w_1^T x = w_2^T x \Rightarrow \|x\|_2\cdot\|W_1\|_2\cdot\cos\theta_1 = \|x\|_2\cdot\|W_2\|_2\cdot\cos\theta_2$

$\|()\|_2$: $L_2$ norm, $0\le\theta_1,\theta_2\le\frac{\pi}{2}$

Decision space:



$w_1$
sample (decision boundary)
$\theta_1$
$\theta_2$
$w_2$

↓

$\|W_1\|_2 \gg \|W_2\|_2$ [long-tail scenarios: weight vectors extremely unbalance]
$\theta_1 \gg \theta_2$

Obviously! decision space for class 2 is pathologically compressed!
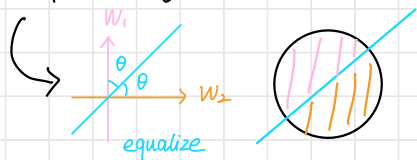
$w_1$
$\theta_1$
$\theta_2$
$w_2$

② Solution

(1) Directly equalize decision space:

↳ Ignore → weight vectors of norm

↳ $L = -\log\left(\dfrac{e^{s\cdot\cos\theta_i}}{\sum_{j=1}^{C} e^{s\cdot\cos\theta_j}}\right)$, $\cos\theta_i = \dfrac{W_i^T x}{\|W_i\|_2 \cdot \|x\|_2}$

↳ optimization goal: Minimum $\theta$ between $W_i$, $X$

↳

equalize



Information Amount: Cat > plane

perceptual manifold volume: cat > plane

model learning focus: cat > plane

So Directly equalize is unreasonable!

- - - - - - - - - - - - - - -

(2) Dynamically adjust decision boundary:

↳ Guided by Category information amount

↳ $L = -\log\left(\dfrac{e^{s\cdot\cos\theta_i}}{e^{s\cdot\cos\theta_i} + \sum_{j=1,\,j\neq i}^{C} e^{s\cdot\cos(\theta_j + m_{ij})}}\right)$

$m_{ij} = \text{Max}\left(0, \frac{1}{\pi}\cdot\log\left(\frac{I_i'}{I_j'}\right)\right)$, $I_i' = \dfrac{e^{e^{\mathcal{I}_i/(\bar{\mathcal{I}}\cdot\sqrt{c})}}}{\sum_{j=1}^{C} e^{\mathcal{I}_i/(\bar{\mathcal{I}}\cdot\sqrt{c})}}\cdot c + 1$, $\bar{\mathcal{I}} = \sum_{i=1}^{c}\mathcal{I}_i$

$I_i'$: normalized information amount of Class $i$

$m_{ij}$: Ratio of information amount of class $i$ and class $j$

↳ $I_i' > I_j'$  decision space expanded

$I_i' < I_j'$  decision space compressed

↳ IGAM Loss making model focus on

complex classes, allocate more decision space.

3: End to end train frame

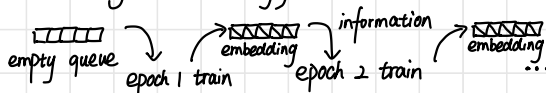<span style="color:red">(1) Dynamic (2) Low-Cost</span>

Engineering challenge: Dynamic update
category information amount → All instances
embedding Covariance matrix

feature slow shift:
sample $a$ in epoch $n$ embedding $\approx$
sample $a$ in epoch $n+1$ embedding

① original strategy:


empty queue    epoch 1 train    epoch 2 train    ...
              embedding    information    embedding

<span style="color:red">defect: storage space ↑</span>

② optimization strategy:

Information Amount ⇐ Global Cov Matrix


Global
local 1    local 2    local 3

C Categories、N instances、d queue length ($d < N$)

Batch → [empty] → [d] → local Cov Matrix, mean

New Batch → [empty] → [  ] → ～

⋮

⇓

<span style="color:red">Local:</span>
Floor $\left(\frac{N}{d}\right) + 1$ local cov matrices $\Sigma_i^k$、local means $\mu_i^k$
$i = 1, \cdots, c$ , $k = 1, \cdots, \text{floor}\left(\frac{N}{d}\right) + 1$

⇓

<span style="color:red">global:</span>
$M_i = \dfrac{1}{N_i}\sum_{k=1}^{\lfloor N/d\rfloor+1} n_i^k \mu_i^k$,    $N_i$: total instances number
                                                                      $n_i^k$: local instances number

$\Sigma_i = \dfrac{1}{N_i}\left(\sum_{k=1}^{\lfloor N/d\rfloor+1} n_i^k \Sigma_i^k + \sum_{k=1}^{\lfloor N/d\rfloor+1} n_i^k (M_i^k - \mu_i)(M_i^k - \mu_i)^T\right)$

⇒ $Vol_i = \frac{1}{2}\log_2 \det(I + \Sigma_i)$