

Semantic Transfer from Head to Tail: Enlarging Tail Margin for Long-Tailed Visual Recognition (Supplementary Materials)

Shan Zhang ^{*,1}, Yao Ni ^{*,1}, Jinhao Du ², Yanxia Liu ^{†,3}, Piotr Koniusz ^{†,4,1}

¹Australian National University, ²Peking University, ³Beijing Union University, ⁴Data61♥CSIRO

¹firstname.lastname@anu.edu.au, ²dujinhao02@gmail.com, ³yanxia.liu@163.com

A. Sketch Proofs in Section 3

A.1. Proof of Assumption 1

Assumption 1. Assume that after augmentation, the feature $\tilde{\mathbf{a}}_i^t$ passed to the classifier, which comes from a tail sample \mathbf{x}_i^t , can be approximately represented by a distribution $\tilde{\mathbf{a}}_i^t \sim \mathcal{N}(\mathbf{a}_i^t, \Delta \Sigma_{th}^i)$. Here, \mathbf{a}_i^t is the feature obtained without augmentation, and $\Delta \Sigma_{th}^i$ is a positive definite covariance matrix.

Proof: For a tail class k_t , each unaugmented feature \mathbf{a}_i^t comes from $\mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$, where $\boldsymbol{\mu}_t$ and $\boldsymbol{\Sigma}_t$ represent the mean and covariance for class k_t . The augmentation operation in Eq. 6 transforms the \mathbf{a}^t into $\tilde{\mathbf{a}}^t$, which follows $\mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_{th})$, with $\boldsymbol{\Sigma}_{th}$ being an approximation to the covariance matrix $\boldsymbol{\Sigma}_h$ (Referring to §A.4, the covariance of augmented tail samples will be close to that of semantically similar head samples under the optimal transformation matrix.). For each i , there is always a suitable matrix $\Delta \Sigma_{th}^i$ that ensures augmented features for class k_t align closely with $\boldsymbol{\Sigma}_{th}$ or $\boldsymbol{\Sigma}_h$. This alignment arises from the synergy of the augmentation method and variations during training, effectively shaping \mathbf{a}_i^t into $\tilde{\mathbf{a}}_i^t \sim \mathcal{N}(\mathbf{a}_i^t, \Delta \Sigma_{th}^i)$.

Remark: Assumption 1 posits that under mild assumptions the augmentation naturally shifts \mathbf{a}_i^t to $\tilde{\mathbf{a}}_i^t \sim \mathcal{N}(\mathbf{a}_i^t, \Delta \Sigma_{th}^i)$.

A.2. Proof of Lemma 2

Lemma 2. Given the negative log softmax function, the loss L_k for samples of class k without augmentation can be derived as:

$$\begin{aligned} L_k &= \frac{1}{n_k} \sum_{i=1}^{n_k} -\log \frac{e^{\mathbf{w}_k^T \mathbf{a}_i + b_k}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{a}_i + b_j}} \\ &= \frac{1}{n_k} \sum_{i=1}^{n_k} \log \left(1 + \sum_{j \neq k} e^{d_i \|\mathbf{w}_j - \mathbf{w}_k\|_2 \cdot \text{sign}(\cos \theta_{i,jk})} \right) \end{aligned} \quad (12)$$

Drawing upon [37], the decision boundary between class j and class k can be formulated as: $(\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{a} + (b_j - b_k) = 0$. d_i is the distance from point \mathbf{a}_i to the decision boundary, $\theta_{i,jk}$ denotes the angle between $\mathbf{w}_j - \mathbf{w}_k$ and \mathbf{a}_i .

^{*}Equal contribution. [†]Corresponding author.

Proof:

$$\begin{aligned}
L_k &= \frac{1}{n_k} \sum_{i=1}^{n_k} -\log \frac{e^{\mathbf{w}_k^T \mathbf{a}_i + b_k}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{a}_i + b_j}} \\
&= \frac{1}{n_k} \sum_{i=1}^{n_k} \log \frac{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{a}_i + b_j}}{e^{\mathbf{w}_k^T \mathbf{a}_i + b_k}} \\
&= \frac{1}{n_k} \sum_{i=1}^{n_k} \log \left(1 + \sum_{j \neq k} e^{(\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{a}_i + (b_j - b_k)} \right) \\
&= \frac{1}{n_k} \sum_{i=1}^{n_k} \log \left(1 + \sum_{j \neq k} e^{d_i \|\mathbf{w}_j - \mathbf{w}_k\|_2 \cdot \text{sign}(\cos \theta_{i,jk})} \right)
\end{aligned} \tag{13}$$

To derive Eq. 13, let us inspect the geometric representation. Focus on the distance d_i from \mathbf{a}_i to the decision boundary $(\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{a}_i + (b_j - b_k) = 0$:

$$d_i = \frac{|(\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{a}_i + (b_j - b_k)|}{\|\mathbf{w}_j - \mathbf{w}_k\|_2} \tag{14}$$

For orientation, when \mathbf{a}_i and \mathbf{w}_k lie on the same side of the decision boundary, the sign of cosine of the angle between $(\mathbf{w}_j - \mathbf{w}_k)$ and \mathbf{a}_i is given by $\text{sign}(\cos \theta_{i,jk}) = -1$. In the opposite scenario, the sign would be positive. Integrating this insight with Eq. 14, we deduce:

$$(\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{a}_i + (b_j - b_k) = d_i \|\mathbf{w}_j - \mathbf{w}_k\|_2 \cdot \text{sign}(\cos \theta_{i,jk}) \tag{15}$$

This directly gives rise to Eq. 13.

A.3. Proof of Theorem 1

Theorem 1. Assume Assumption 1 holds when using our augmentation. The loss function L_k^t for tail class k is:

$$\begin{aligned}
L_k^t &= \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbb{E}_{\tilde{\mathbf{a}}_i^t} \left[-\log \frac{e^{\mathbf{w}_k^T \tilde{\mathbf{a}}_i^t + b_k}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \tilde{\mathbf{a}}_i^t + b_j}} \right] \\
&\leq \frac{1}{n_k} \sum_{i=1}^{n_k} \log \left(1 + \sum_{j \neq k} \beta_{jk}^i e^{d_i \|\mathbf{w}_j - \mathbf{w}_k\|_2 \cdot \text{sign}(\cos \theta_{i,jk})} \right)
\end{aligned} \tag{16}$$

where $\beta_{jk}^i = e^{\frac{1}{2}(\mathbf{w}_j - \mathbf{w}_k)^T \Delta \Sigma_{th}^i (\mathbf{w}_j - \mathbf{w}_k)}$. Furthermore,

$$\beta_{jk}^i = \exp \left(\frac{1}{2} \mathbf{v}_{jk}^i{}^T \mathbf{\Lambda}^i \mathbf{v}_{jk}^i \right) > 1 \tag{17}$$

where $\mathbf{V}^i \mathbf{\Lambda}^i \mathbf{V}^i{}^T = \Delta \Sigma_{th}^i$ and $\mathbf{v}_{jk}^i = \mathbf{V}^i{}^T (\mathbf{w}_j - \mathbf{w}_k)$.

Proof of Eq. 16:

$$\begin{aligned}
L_k^t &= \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbb{E}_{\tilde{\mathbf{a}}_i^t} \left[-\log \frac{e^{\mathbf{w}_k^T \tilde{\mathbf{a}}_i^t + b_k}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \tilde{\mathbf{a}}_i^t + b_j}} \right] \\
&= \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbb{E}_{\tilde{\mathbf{a}}_i^t} \left[\log \left(1 + \sum_{j \neq k} e^{(\mathbf{w}_j - \mathbf{w}_k)^T \tilde{\mathbf{a}}_i^t + (b_j - b_k)} \right) \right] \\
&\leq \frac{1}{n_k} \sum_{i=1}^{n_k} \log \left(1 + \sum_{j \neq k} \mathbb{E}_{\tilde{\mathbf{a}}_i^t} \left[e^{(\mathbf{w}_j - \mathbf{w}_k)^T \tilde{\mathbf{a}}_i^t + (b_j - b_k)} \right] \right) \tag{18}
\end{aligned}$$

$$= \frac{1}{n_k} \sum_{i=1}^{n_k} \log \left(1 + \sum_{j \neq k} e^{(\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{a}_i^t + (b_j - b_k) + \frac{1}{2} (\mathbf{w}_j - \mathbf{w}_k)^T \Delta \Sigma_{th} (\mathbf{w}_j - \mathbf{w}_k)} \right) \tag{19}$$

$$\begin{aligned}
&= \frac{1}{n_k} \sum_{i=1}^{n_k} \log \left(1 + \sum_{j \neq k} e^{\frac{1}{2} (\mathbf{w}_j - \mathbf{w}_k)^T \Delta \Sigma_{th}^i (\mathbf{w}_j - \mathbf{w}_k)} \cdot e^{(\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{a}_i^t + (b_j - b_k)} \right) \\
&= \frac{1}{n_k} \sum_{i=1}^{n_k} \log \left(1 + \sum_{j \neq k} \beta_{jk} e^{(\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{a}_i^t + (b_j - b_k)} \right) \\
&= \frac{1}{n_k} \sum_{i=1}^{n_k} \log \left(1 + \sum_{j \neq k} \beta_{jk} e^{d_i \|\mathbf{w}_j - \mathbf{w}_k\|_2 \cdot \text{sign}(\cos \theta_{i,k,j})} \right). \tag{20}
\end{aligned}$$

In the above derivation, the inequality Eq. 18 is a direct consequence of Jensen's inequality $\mathbb{E}[\log X] \leq \log \mathbb{E}[X]$. Eq. 19 is obtained by leveraging the moment-generating function $\mathbb{E}[e^{tX}] = e^{t\mu + \frac{1}{2}\sigma^2 t^2}$ where $X \sim \mathcal{N}(\mu, \sigma^2)$, and the fact that $(\mathbf{w}_j - \mathbf{w}_k)^T \tilde{\mathbf{a}}_i^t + (b_j - b_k)$ is a Gaussian random variable drawn from $\mathcal{N}\left((\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{a}_i^t + (b_j - b_k), (\mathbf{w}_j - \mathbf{w}_k)^T \Delta \Sigma_{th}^i (\mathbf{w}_j - \mathbf{w}_k)\right)$. Lastly, Eq. 20 is derived by incorporating Eq. 15.

Proof of Eq. 17: Performing SVD on the positive definite symmetric covariance matrix $\Delta \Sigma_{th}^i$, we obtain $\Delta \Sigma_{th}^i = \mathbf{V}^i \mathbf{\Lambda}^i \mathbf{V}^{iT}$, where \mathbf{V}^i represents the eigenvectors and $\mathbf{\Lambda}^i$ is the diagonal matrix of eigenvalues. By incorporating $\mathbf{w}_j - \mathbf{w}_k$ and \mathbf{V}^i into a single term, we define $\mathbf{v}_{jk}^i = \mathbf{V}^{iT} (\mathbf{w}_j - \mathbf{w}_k)$. On deriving β_{jk}^i , we get:

$$\begin{aligned}
\beta_{jk}^i &= e^{\frac{1}{2} (\mathbf{w}_j - \mathbf{w}_k)^T \Delta \Sigma_{th}^i (\mathbf{w}_j - \mathbf{w}_k)} \\
&= e^{\frac{1}{2} (\mathbf{w}_j - \mathbf{w}_k)^T \mathbf{V}^i \mathbf{\Lambda}^i \mathbf{V}^{iT} (\mathbf{w}_j - \mathbf{w}_k)} \\
&= e^{\frac{1}{2} (\mathbf{V}^{iT} (\mathbf{w}_j - \mathbf{w}_k))^T \mathbf{\Lambda}^i (\mathbf{V}^{iT} (\mathbf{w}_j - \mathbf{w}_k))} \\
&= e^{\frac{1}{2} \mathbf{v}_{jk}^{iT} \mathbf{\Lambda}^i \mathbf{v}_{jk}^i}
\end{aligned}$$

Let v_c be the c^{th} element of \mathbf{v}_{jk}^i and $\lambda_c \geq 0$ be the c^{th} $\text{diag}(\mathbf{\Lambda}^i)$, we have:

$$\beta_{jk}^i = e^{\frac{1}{2} \sum_c v_c^2 \lambda_c} > e^0.$$

For any non-zero vector \mathbf{w} and a positive definite matrix \mathbf{A} , the result $\mathbf{w}^T \mathbf{A} \mathbf{w} > 0$ implies $\beta_{jk}^i > 1$ due to the non-zero vector $\mathbf{w}_j - \mathbf{w}_k$ and positive definiteness of $\Delta \Sigma_{th}^i$. Furthermore, the relationship $\sum_c \lambda_c = \text{trace}(\mathbf{\Lambda}^i) = \text{trace}(\Delta \Sigma_{th}^i)$ indicates the larger the semantic similarities between the head and tail samples, and the more diverse the head class, the greater the value of β_{jk}^i .

A.4. Design of transformation matrix

We aim to design a transformation matrix such that the covariance of transformed tail feature \mathbf{F}_t , aligns closely with the covariance of head samples. The objective is formulated as:

$$\tilde{\mathbf{F}}_t^* = \arg \min_{\tilde{\mathbf{F}}_t} \|\tilde{\mathbf{F}}_t^T \tilde{\mathbf{F}}_t - \mathbf{F}_h^T \mathbf{F}_h\|_F^2 \tag{21}$$

$$s.t. \quad \tilde{\mathbf{F}}_t = \mathbf{T} \mathbf{F}_t. \tag{22}$$

Substituting the constraint from Eq. 22 into Eq. 21, we find optimality at:

$$\mathbf{F}_t^T \mathbf{T}^T \mathbf{T} \mathbf{F}_t = \mathbf{F}_h^T \mathbf{F}_h. \quad (23)$$

Upon applying singular value decomposition (SVD) to \mathbf{F}_t and \mathbf{F}_h , yielding $\mathbf{V}_t \boldsymbol{\Sigma}_t \mathbf{V}_t^T$ and $\mathbf{V}_h \boldsymbol{\Sigma}_h \mathbf{V}_h^T$ and insert them into Eq. 23, a solution set emerges:

$$\mathbf{T} = (\mathbf{V}_h \boldsymbol{\Sigma}_h^{\frac{1}{2}} \mathbf{V}_h^T) \mathbf{U} (\mathbf{V}_t \boldsymbol{\Sigma}_t^{-\frac{1}{2}} \mathbf{V}_t^T)^T, \quad (24)$$

where $\mathbf{U} \in \mathbb{R}^{C \times C}$ is a orthogonal group. Eq. 24 indicates that the transformation matrix, \mathbf{T} , is influenced by the covariance matrices of both tail and head classes. Empirical validation of this design is presented in Sec. 4.3.