

# Supplementary Material for Rotation Regularization Without Rotation

Takumi Kobayashi

takumi.kobayashi@aist.go.jp

## 1 Proofs

We denote vectors by bold lowercase letters, e.g.,  $\mathbf{x}$ , and normalized vectors by using  $\bar{\cdot}$ , e.g.,  $\bar{\mathbf{x}} = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}$ ; thus,  $\mathbf{x} \in \mathbb{R}^D \Rightarrow \bar{\mathbf{x}} \in \mathbb{S}^{D-1}$ , the hyper-sphere in  $\mathbb{R}^D$ . Matrices are denoted by bold uppercase letters, e.g.,  $\mathbf{R}$ .

**Lemma 1.** *A vector  $\bar{\mathbf{x}} \in \mathbb{S}^{D-1}$  is rotated by an angle  $\alpha$  through a rotation matrix  $\mathbf{R}_\alpha$ . So rotated vector is described by using a differential vector  $\exists \bar{\mathbf{z}} \in \mathbb{S}^{D-2}$  which is in the orthogonal complement space to the input vector  $\bar{\mathbf{x}}$  as*

$$\mathbf{R}_\alpha \bar{\mathbf{x}} = \cos \alpha \bar{\mathbf{x}} + \sin \alpha \bar{\mathbf{z}}, \text{ where } \|\bar{\mathbf{z}}\|_2 = 1, \bar{\mathbf{x}}^\top \bar{\mathbf{z}} = 0. \quad (1)$$

*Proof.* Fig. 1 would be helpful to grasp relationships among the following vectors. Let the rotated vector be denoted by  $\bar{\boldsymbol{\xi}} = \mathbf{R}_\alpha \bar{\mathbf{x}}$  which satisfies

$$\bar{\mathbf{x}}^\top \bar{\boldsymbol{\xi}} = \cos \alpha. \quad (2)$$

We define the differential vector  $\mathbf{z}$  as

$$\mathbf{z} = \bar{\boldsymbol{\xi}} - \cos \alpha \bar{\mathbf{x}}, \quad (3)$$

and it has the following properties;

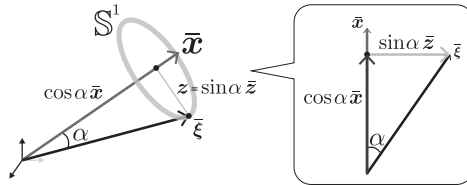
$$\bar{\mathbf{x}}^\top \mathbf{z} = \bar{\mathbf{x}}^\top \bar{\boldsymbol{\xi}} - \cos \alpha \|\bar{\mathbf{x}}\|_2^2 = 0 \quad (4)$$

$$\|\mathbf{z}\|_2^2 = \|\bar{\boldsymbol{\xi}}\|_2^2 + \cos^2 \alpha \|\bar{\mathbf{x}}\|_2^2 - 2 \cos \alpha \bar{\mathbf{x}}^\top \bar{\boldsymbol{\xi}} = 1 - \cos^2 \alpha = \sin^2 \alpha, \quad (5)$$

where we apply (2). Since  $0 \leq \alpha < \frac{\pi}{2}$ , the differential vector  $\mathbf{z}$  is thus described by using  $\bar{\mathbf{z}} = \frac{\mathbf{z}}{\|\mathbf{z}\|_2}$  as

$$\mathbf{z} = \sin \alpha \bar{\mathbf{z}} \text{ where } \|\bar{\mathbf{z}}\|_2 = 1, \bar{\mathbf{x}}^\top \bar{\mathbf{z}} = 0. \quad (6)$$

It should be noted that  $\bar{\mathbf{z}}$  is laid on the hyper-sphere  $\mathbb{S}^{D-2}$  embedded in the  $D - 1$  dimensional subspace orthogonal to  $\bar{\mathbf{x}} \in \mathbb{R}^D$ .  $\square$



**Fig. 1.** Reparameterization of a rotated vector.

**Lemma 2.** *Projection of random vectors  $\bar{\mathbf{a}}$  uniformly distributed on a unit hypersphere  $\mathbb{S}^{m-1}$  into a unit-length vector  $\bar{\mathbf{b}} \in \mathbb{S}^{m-1}$  follows Beta distribution.*

$$\bar{\mathbf{a}} \sim \text{Unif}(\mathbb{S}^{m-1}), \bar{\mathbf{b}} \in \mathbb{S}^{m-1}, u = \frac{1 + \bar{\mathbf{a}}^\top \bar{\mathbf{b}}}{2} \Rightarrow u \sim \text{Beta}\left(\frac{m-1}{2}, \frac{m-1}{2}\right). \quad (7)$$

For the higher dimensional case  $m \gg 1$ , it approaches Gaussian distribution as

$$\bar{\mathbf{a}}^\top \bar{\mathbf{b}} \sim \mathcal{N}\left(0, \frac{1}{\sqrt{m}}\right) \quad (8)$$

*Proof.* As shown in Fig. 2, we consider the probability density of  $dt$  region at  $t = \bar{\mathbf{a}}^\top \bar{\mathbf{b}}$  along  $\bar{\mathbf{b}}$ . It is proportional to a surface volume  $dV$  on  $\mathbb{S}^{m-1}$ , a gray-colored belt in Fig. 2 which is composed of the length  $\frac{dt}{\sqrt{1-t^2}}$  and the hyper volume of  $\mathbb{S}^{m-2}$  with radius  $\sqrt{1-t^2}$  and thereby computed as

$$\mathbf{p}(t)dt \propto dV \propto (\sqrt{1-t^2})^{m-2} \frac{dt}{\sqrt{1-t^2}} = (1-t^2)^{\frac{m-3}{2}} dt. \quad (9)$$

Due to  $\bar{\mathbf{a}}^\top \bar{\mathbf{b}} = t = 2u - 1$ , the probability density function  $\mathbf{q}(u)$  is described by

$$\mathbf{q}(u)du = \mathbf{p}(t)dt \propto \{1 - (2u - 1)^2\}^{\frac{m-3}{2}} \cdot 2du = 2(4u - 4u^2)^{\frac{m-3}{2}} du \quad (10)$$

$$= 2^{m-2} u^{\frac{m-3}{2}} (1-u)^{\frac{m-3}{2}} du \propto u^{\frac{m-1}{2}-1} (1-u)^{\frac{m-1}{2}-1} du, \quad (11)$$

which corresponds to the Beta distribution,  $\text{Beta}(u; \beta, \beta) \propto u^{\beta-1} (1-u)^{\beta-1}$  with  $\beta = \frac{m-1}{2}$ . Thus, the mean and variance of  $u$  are

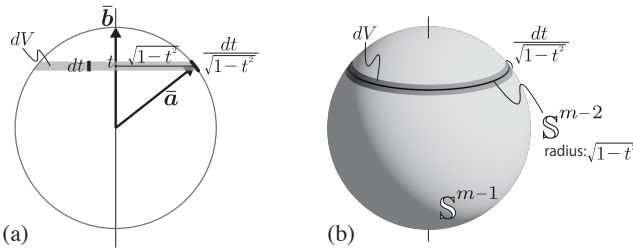
$$\mathbf{E}_{u \sim \text{Beta}}[u] = \frac{1}{2}, \quad \text{Var}_{u \sim \text{Beta}}[u] = \frac{\beta^2}{4\beta^2(2\beta+1)} = \frac{1}{4(2\beta+1)} = \frac{1}{4m}. \quad (12)$$

Asymptotic property of the Beta distribution as  $\beta \rightarrow \infty$  is

$$\text{Beta}(\beta, \beta) \rightarrow \mathcal{N}\left(\frac{1}{2}, \frac{1}{2\sqrt{2\beta+1}}\right), \quad (13)$$

the proof of which is found such as in [7]. It is applied to  $u \sim \text{Beta}(\beta, \beta)$  with  $\beta = \frac{m-1}{2}$  to produce

$$\bar{\mathbf{a}}^\top \bar{\mathbf{b}} = 2u - 1 \sim \mathcal{N}\left(0, \frac{1}{\sqrt{m}}\right). \quad \square \quad (14)$$



**Fig. 2.** Projection from a hyper sphere  $\mathbb{S}^{m-1}$  into a unit-length vector  $\bar{\mathbf{b}}$ .

**Theorem 1.** Random rotation matrix  $\mathbf{R}_\alpha$  of an angle  $\alpha$  is applied to an inner product between two unit-length vectors  $\bar{\mathbf{w}} \in \mathbb{R}^D$  and  $\bar{\mathbf{x}} \in \mathbb{R}^D$  where  $\bar{\mathbf{w}}^\top \bar{\mathbf{x}} = \cos \theta$ . Then, the inner product is endowed with stochasticity by the random  $\mathbf{R}_\alpha$  and is statistically described by

$$\bar{\mathbf{w}}^\top \mathbf{R}_\alpha \bar{\mathbf{x}} = \cos \alpha \cos \theta + (2\eta - 1) \sin \alpha \sin \theta \text{ where } \eta \sim \text{Beta}\left(\frac{D-2}{2}, \frac{D-2}{2}\right). \quad (15)$$

For the higher dimensional case  $D \gg 1$ , it approaches Gaussian distribution as

$$\bar{\mathbf{w}}^\top \mathbf{R}_\alpha \bar{\mathbf{x}} = \cos \alpha \cos \theta + \frac{\epsilon}{\sqrt{D-1}} \sin \alpha \sin \theta \text{ where } \epsilon \sim \mathcal{N}(0, 1). \quad (16)$$

*Proof.* By applying Lemma 1, we can obtain

$$\bar{\mathbf{w}}^\top \mathbf{R}_\alpha \bar{\mathbf{x}} = \bar{\mathbf{w}}^\top (\cos \alpha \bar{\mathbf{x}} + \sin \alpha \bar{\mathbf{z}}) = \cos \alpha \cos \theta + \sin \alpha \bar{\mathbf{w}}^\top \bar{\mathbf{z}}. \quad (17)$$

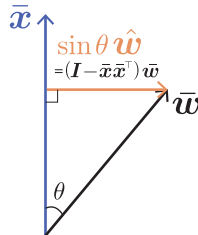
Due to the randomness of the rotation matrix  $\mathbf{R}_\alpha$ ,  $\bar{\mathbf{z}}$  is uniformly drawn from  $\mathbb{S}^{D-2}$  embedded in the subspace  $\mathbb{R}^{D-1}$  perpendicular to  $\bar{\mathbf{x}}$ ; so,  $\bar{\mathbf{z}} = (\mathbf{I} - \bar{\mathbf{x}}\bar{\mathbf{x}}^\top)\bar{\mathbf{z}}$ . In that subspace, the inner product  $\bar{\mathbf{w}}^\top \bar{\mathbf{z}}$  is described by

$$\bar{\mathbf{w}}^\top \bar{\mathbf{z}} = \bar{\mathbf{w}}^\top (\mathbf{I} - \bar{\mathbf{x}}\bar{\mathbf{x}}^\top)\bar{\mathbf{z}} = \sin \theta \hat{\mathbf{w}}^\top \bar{\mathbf{z}}, \quad (18)$$

where  $\hat{\mathbf{w}} = \frac{(\mathbf{I} - \bar{\mathbf{x}}\bar{\mathbf{x}}^\top)\bar{\mathbf{w}}}{\|(\mathbf{I} - \bar{\mathbf{x}}\bar{\mathbf{x}}^\top)\bar{\mathbf{w}}\|_2}$  is in the subspace  $\mathbb{R}^{D-1}$  and  $\|(\mathbf{I} - \bar{\mathbf{x}}\bar{\mathbf{x}}^\top)\bar{\mathbf{w}}\|_2 = \sin \theta$  as shown in Fig. 3. Then, based on the stochasticity  $\bar{\mathbf{z}} \sim \text{Unif}(\mathbb{S}^{D-2})$ , Lemma 2 with  $m = D - 1$  statistically rewrites the inner product into

$$\bar{\mathbf{w}}^\top \bar{\mathbf{z}} = \sin \theta (2\eta - 1), \quad (19)$$

where  $\eta \sim \text{Beta}(\frac{D-2}{2}, \frac{D-2}{2})$ . (17) and (19) lead to (15) and the asymptotic property of Beta distribution shown in Lemma 2 produces (16) with  $m = D - 1$ .  $\square$



**Fig. 3.** Projection from  $\bar{\mathbf{w}}$  into the subspace perpendicular to  $\bar{\mathbf{x}}$ .

## 2 Hyper parameters of comparison methods

In DropOut [13], we employ small probability  $p = 0.2$  to mask feature elements according to the analysis [8] and preliminary experiments. For the margin-based losses, based on the preliminary experiments, we apply angular margin parameter of  $m = 0.1$  in ArcFace [2],  $m = 0.1$  in CosFace [14] and  $\alpha = \sqrt{0.1}$  in NoisySoftmax [1].

## 3 Datasets

The details of the datasets that we use in Sec. 5 of the main manuscript are shown below. Except for CALTECH101 [3], we use the train/test splits provided in the respective datasets; for SUN397 [16], we use the first split out of 10 splits given in the dataset. In CALTECH101, following the standard protocol, we randomly draw 30 training samples per category and use the remaining samples as test.

(a) Long-tailed datasets					
	IMAGENET-LT [9]	iNAT2018 [5]	PLACES-LT [9]		
category	1000 objects	8142 species	365 scenes		
training sample	115,846	437,513	62,500		
test sample	50,000	24,426	36,500		
majority : minority	1,280 : 5	1,000 : 2	4,980 : 5		

(b) Downstream datasets					
	CUB200 [15]	AIRCRAFT100 [11]	CAR196 [6]	SUN397 [16]	CALTECH101 [3]
category	200 birds	100 planes	196 cars	397 scenes	101 objects
trn. sample	5,994	6,667	8,144	19,850	3,030
tst. sample	5,794	3,333	8,041	19,850	5,647

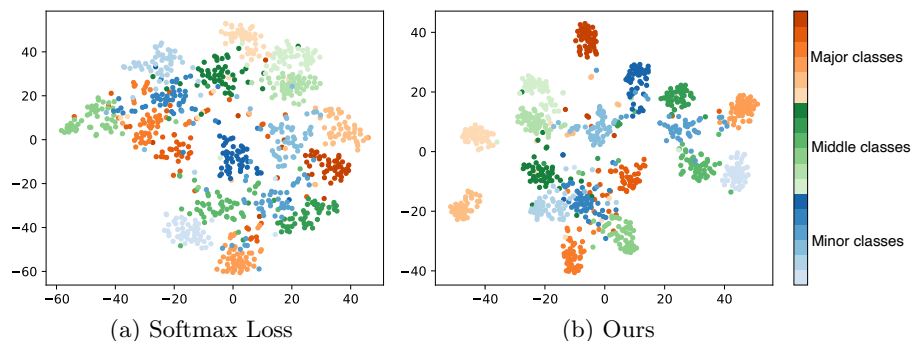
  

(c) Person re-identification datasets		
	MARKET1501 [17]	DUKEMTMC [12]
trn. category	750 identities	702 identities
trn. sample	12,936	16,522
tst. category	751 identities	702 identities
query sample	3,368	2,228
gallery sample	19,732	17,661

## 4 Discriminativity of feature representation

We apply t-SNE [10] to show the feature distributions of ResNet10 backbone trained on IMAGENET-LT dataset by using the baseline softmax loss and ours with statistical rotation regularization. Fig. 4 demonstrates that our method improves feature distribution in comparison to the softmax loss.

The discriminativity of features is quantitatively measured by means of discriminant score [4], the ratio of between-class feature variance to within-class one;  $\text{tr}(\Sigma_B)/\text{tr}(\Sigma_W)$ . Table 1 demonstrates that our method improves the score, contributing to intra-class compactness as well as inter-class separability.



**Fig. 4.** Visualization of ResNet10 feature distributions on IMAGENET-LT via t-SNE [10]. Each point indicates a sample drawn from major, middle and minor classes on the validation set.

**Table 1.** Discriminant score  $\text{tr}(\Sigma_B)/\text{tr}(\Sigma_W)$ . Higher is better.

	IMAGENET-LT	<i>i</i> NAT2018	PLACES-LT	
	ResNet10	ResNet50	ResNet10	ResNet152
Softmax Loss	0.385	1.475	0.270	0.372
Ours	<b>0.692</b>	<b>1.747</b>	<b>0.495</b>	<b>0.672</b>

## References

1. Chen, B., Deng, W., Du, J.: Noisy softmax: Improving the generalization ability of dcnn via postponing the early softmax saturation. In: CVPR. pp. 4021–4030 (2017)
2. Deng, J., Guo, J., Niannan, X., Zafeiriou, S.: Arcface: Additive angular margin loss for deep face recognition. In: CVPR (2019)

3. Fei-Fei, L., Fergus, R., Perona, P.: Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In: Computer Vision and Pattern Recognition Workshop (2004)
4. Fukunaga, K.: Introduction to Statistical Pattern Recognition. Academic Press, Boston (1990)
5. iNaturalist: The inaturalist 2018 competition dataset. [https://github.com/visipedia/inat\\_comp/tree/master/2018](https://github.com/visipedia/inat_comp/tree/master/2018) (2018)
6. Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3d object representations for fine-grained categorization. In: Workshop on 3D Representation and Recognition (2013)
7. Leemis, L.: <http://www.math.wm.edu/~leemis/chart/UDR/PDFs/BetaNormal.pdf>
8. Li, X., Chen, S., Hu, X., Yang, J.: Understanding the disharmony between dropout and batch normalization by variance shift. In: CVPR. pp. 2682–2690 (2019)
9. Liu, Z., Miao, Z., Zhan, X., Wang, J., Gong, B., Yu, S.X.: Large-scale long-tailed recognition in an open world. In: CVPR. pp. 2537–2546 (2019)
10. van der Maaten, L., Hinton, G.: Visualizing data using t-SNE. *Journal of Machine Learning Research* **9**, 2579–2605 (2008)
11. Maji, S., Rahtu, E., Kannala, J., Blaschko, M.B., Vedaldi, A.: Fine-grained visual classification of aircraft. arXiv:1306.5151 (2013)
12. Ristani, E., Solera, F., Zou, R., Cucchiara, R., Tomasi, C.: Performance measures and a data set for multi-target, multi-camera tracking. In: ECCV Workshop (2016)
13. Srivastava, N., Hinton, G.E., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout : A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* **15**, 1929–1958 (2014)
14. Wang, H., Wang, Y., Zhou, Z., Ji, X., Gong, D., Zhou, J., Li, Z., Liu, W.: Cosface: Large margin cosine loss for deep face recognition. In: CVPR. pp. 5265–5274 (2018)
15. Welinder, P., Branson, S., Mita, T., Wah, C., Schroff, F., Belongie, S., Perona, P.: Caltech-UCSD Birds 200. Tech. Rep. CNS-TR-2010-001, California Institute of Technology (2010)
16. Xiao, J., Hays, J., Ehinger, K.A., Oliva, A., Torralba, A.: Sun database: Large-scale scene recognition from abbey to zoo. In: CVPR (2010)
17. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: ICCV (2015)