

---

# Supplementary Material for Mean Field Theory in Deep Metric Learning

---

Anonymous Author(s)

Affiliation

Address

email

## A Additional experimental results

In this section, we present the experimental results that cannot be shown in the main paper due to the page limit.

**MLRC results.** Tables 1 – 3 show the complete results of Table 1 in the main paper, which are obtained the modern benchmark protocol proposed in the “Metric Learning Reality Check” (MLRC) paper [1]. In the CUB-200-2011 (CUB) dataset [2], MeanFieldContrastive (MFCont.) and MeanFieldClassWiseMultiSimilarity (MFCWMS) losses outperform the others in Mean Average Precision at R (MAP@R) and R-Precision (RP), while ProxyAnchor loss [3] is better in Precision at 1 (P@1) in the separated case. In contrast, in the Stanford Online Products (SOP) dataset [4], the MFCWMS loss shows the best performance in all the metrics.

**Learning curves.** Figure 1 shows learning curves obtained in the traditional evaluation protocol [3, 5] in fixed seeds, associated with Tables 2 and 3 in the main paper. Both MFCont. and MFCWMS losses show faster convergence than the ProxyAnchor loss. In the smaller datasets (CUB and Cars), accuracies of our mean field losses seem to decrease faster while we don’t see such behaviors in the larger datasets (SOP and InShop [6]). This phenomenon might be caused by strong repulsive interactions with negative mean fields. For larger datasets, the embedding spaces may be sufficiently populated to balance the repulsive force, while this may not be the case for smaller datasets. It might not occur for ProxyAnchor loss since repulsive forces for ProxyAnchor loss are weighted depending on distances between proxy and negative samples.

**Impact of batch size in InShop.** Table 4 compares the MAP@R in ProxyAnchor and MFCWMS losses in the InShop dataset varying the batch size. As mentioned in the main paper, the accuracy of ProxyAnchor starts to decrease gradually for large batch sizes, while that of MFCWMS loss drops at batch size 150. Moreover, Table 5 shows the MAP@R in ProxyAnchor and MFCWMS for the InShop dataset without the query-gallery split of test data. In this case, accuracies of both losses start to decrease gradually around batch size 150. Thus, we conclude the accuracy drop in the MFCWMS loss probably comes from the specific query-gallery split.

Table 1: MLRC evaluation results in CUB-200-2011 [2]. We carry out 10 test runs and show averaged metrics with their confidence intervals.

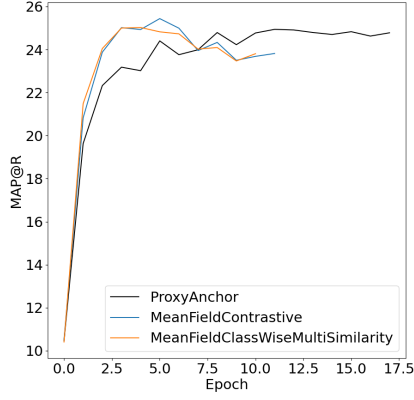
Loss	Separated (128D)			Concatenated (512D)		
	MAP@R	P@1	RP	MAP@R	P@1	RP
ArcFace	21.46 $\pm$ 0.13	59.98 $\pm$ 0.22	32.31 $\pm$ 0.14	26.39 $\pm$ 0.16	67.11 $\pm$ 0.23	37.23 $\pm$ 0.17
CosFace	21.19 $\pm$ 0.22	59.74 $\pm$ 0.28	32.00 $\pm$ 0.23	<b>26.54 <math>\pm</math> 0.29</b>	67.14 $\pm$ 0.29	<b>37.38 <math>\pm</math> 0.28</b>
MS	20.98 $\pm$ 0.16	59.38 $\pm$ 0.27	31.84 $\pm$ 0.15	26.20 $\pm$ 0.16	67.34 $\pm$ 0.35	36.99 $\pm$ 0.16
MS+Miner	20.78 $\pm$ 0.17	59.02 $\pm$ 0.25	31.67 $\pm$ 0.16	25.94 $\pm$ 0.18	67.08 $\pm$ 0.32	36.77 $\pm$ 0.16
ProxyNCA	18.75 $\pm$ 0.18	57.06 $\pm$ 0.27	29.64 $\pm$ 0.21	23.84 $\pm$ 0.22	65.60 $\pm$ 0.28	34.82 $\pm$ 0.25
ProxyAnch.	<b>21.67 <math>\pm</math> 0.22</b>	<b>60.80 <math>\pm</math> 0.33</b>	<b>32.53 <math>\pm</math> 0.23</b>	26.48 $\pm$ 0.23	<b>67.72 <math>\pm</math> 0.30</b>	37.30 $\pm$ 0.23
Cont.	21.02 $\pm$ 0.14	59.35 $\pm$ 0.33	31.80 $\pm$ 0.15	26.37 $\pm$ 0.18	<b>67.67 <math>\pm</math> 0.25</b>	37.10 $\pm$ 0.19
MFCCont.	<b>22.01 <math>\pm</math> 0.10</b>	<b>60.29 <math>\pm</math> 0.23</b>	<b>32.85 <math>\pm</math> 0.10</b>	<b>27.16 <math>\pm</math> 0.07</b>	67.64 $\pm$ 0.27	<b>37.95 <math>\pm</math> 0.07</b>
CWMS	21.48 $\pm$ 0.27	60.09 $\pm$ 0.27	32.32 $\pm$ 0.26	26.94 $\pm$ 0.29	<b>68.24 <math>\pm</math> 0.42</b>	37.69 $\pm$ 0.27
MFCWMS	<b>22.11 <math>\pm</math> 0.08</b>	<b>60.28 <math>\pm</math> 0.10</b>	<b>32.96 <math>\pm</math> 0.08</b>	<b>27.03 <math>\pm</math> 0.12</b>	67.63 $\pm$ 0.21	<b>37.83 <math>\pm</math> 0.12</b>

Table 2: MLRC evaluation results in Cars-196 [7]. We carry out 10 test runs and show averaged metrics with their confidence intervals.

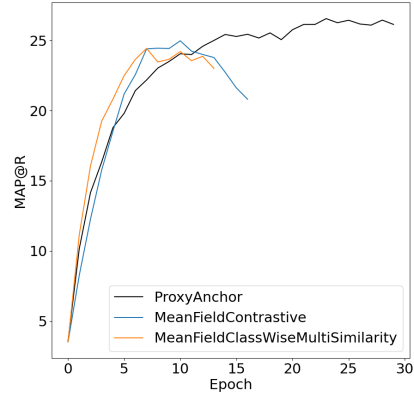
Loss	Separated (128D)			Concatenated (512D)		
	MAP@R	P@1	RP	MAP@R	P@1	RP
ArcFace	18.25 $\pm$ 0.12	71.12 $\pm$ 0.36	28.63 $\pm$ 0.13	<b>27.63 <math>\pm</math> 0.15</b>	84.39 $\pm$ 0.15	<b>37.45 <math>\pm</math> 0.15</b>
CosFace	18.49 $\pm$ 0.13	74.66 $\pm$ 0.21	28.75 $\pm$ 0.12	26.96 $\pm$ 0.25	85.29 $\pm$ 0.26	36.80 $\pm$ 0.24
MS	18.66 $\pm$ 0.30	71.89 $\pm$ 0.33	29.42 $\pm$ 0.29	27.19 $\pm$ 0.41	84.03 $\pm$ 0.30	37.39 $\pm$ 0.36
MS+Miner	18.49 $\pm$ 0.23	71.99 $\pm$ 0.28	29.20 $\pm$ 0.23	26.89 $\pm$ 0.38	83.89 $\pm$ 0.36	37.09 $\pm$ 0.33
ProxyNCA	17.43 $\pm$ 0.11	70.96 $\pm$ 0.26	27.85 $\pm$ 0.10	26.78 $\pm$ 0.18	84.31 $\pm$ 0.24	36.83 $\pm$ 0.17
ProxyAnch.	<b>19.44 <math>\pm</math> 0.17</b>	<b>76.15 <math>\pm</math> 0.25</b>	<b>29.89 <math>\pm</math> 0.18</b>	26.81 $\pm$ 0.27	<b>85.53 <math>\pm</math> 0.30</b>	36.76 $\pm$ 0.26
Cont.	17.04 $\pm$ 0.26	69.77 $\pm$ 0.40	27.48 $\pm$ 0.26	24.93 $\pm$ 0.46	81.87 $\pm$ 0.35	35.12 $\pm$ 0.42
MFCCont.	<b>18.12 <math>\pm</math> 0.13</b>	<b>71.77 <math>\pm</math> 0.28</b>	<b>28.54 <math>\pm</math> 0.14</b>	<b>27.37 <math>\pm</math> 0.18</b>	<b>84.56 <math>\pm</math> 0.21</b>	<b>37.19 <math>\pm</math> 0.18</b>
CWMS	<b>19.27 <math>\pm</math> 0.26</b>	<b>74.19 <math>\pm</math> 0.30</b>	<b>29.95 <math>\pm</math> 0.25</b>	<b>27.80 <math>\pm</math> 0.33</b>	<b>85.18 <math>\pm</math> 0.28</b>	<b>37.89 <math>\pm</math> 0.29</b>
MFCWMS	18.85 $\pm$ 0.16	73.02 $\pm$ 0.20	29.55 $\pm$ 0.15	26.98 $\pm$ 0.31	84.00 $\pm$ 0.22	37.11 $\pm$ 0.27

Table 3: MLRC evaluation results in Stanford Online Products [4]. We carry out 10 test runs and show averaged metrics with their confidence intervals. We remove ProxyAnchor because it fails to converge in our settings.

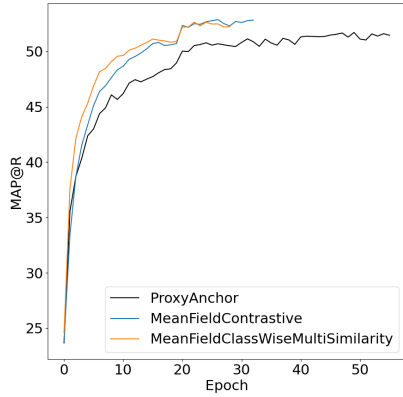
Loss	Separated (128D)			Concatenated (512D)		
	MAP@R	P@1	RP	MAP@R	P@1	RP
ArcFace	41.47 $\pm$ 0.24	71.39 $\pm$ 0.20	44.35 $\pm$ 0.23	<b>47.37 <math>\pm</math> 0.23</b>	<b>76.13 <math>\pm</math> 0.16</b>	<b>50.22 <math>\pm</math> 0.22</b>
CosFace	41.01 $\pm$ 0.24	71.03 $\pm$ 0.22	43.89 $\pm$ 0.24	46.77 $\pm$ 0.20	75.69 $\pm$ 0.13	49.63 $\pm$ 0.20
MS	41.87 $\pm$ 0.21	71.10 $\pm$ 0.18	45.00 $\pm$ 0.20	46.70 $\pm$ 0.18	75.21 $\pm$ 0.15	49.70 $\pm$ 0.17
MS+Miner	41.90 $\pm$ 0.30	71.08 $\pm$ 0.25	45.05 $\pm$ 0.30	46.57 $\pm$ 0.28	75.09 $\pm$ 0.19	49.57 $\pm$ 0.28
ProxyNCA	<b>42.73 <math>\pm</math> 0.11</b>	<b>71.77 <math>\pm</math> 0.08</b>	<b>45.72 <math>\pm</math> 0.11</b>	46.73 $\pm$ 0.13	75.24 $\pm$ 0.10	49.61 $\pm$ 0.13
Cont.	41.09 $\pm$ 0.18	70.04 $\pm$ 0.16	44.18 $\pm$ 0.19	45.35 $\pm$ 0.19	73.88 $\pm$ 0.15	48.28 $\pm$ 0.19
MFCCont.	<b>43.62 <math>\pm</math> 0.36</b>	<b>72.74 <math>\pm</math> 0.29</b>	<b>46.55 <math>\pm</math> 0.35</b>	<b>47.01 <math>\pm</math> 0.21</b>	<b>75.57 <math>\pm</math> 0.16</b>	<b>49.85 <math>\pm</math> 0.20</b>
CWMS	41.53 $\pm$ 0.20	70.76 $\pm$ 0.16	44.50 $\pm$ 0.21	45.13 $\pm$ 0.16	73.99 $\pm$ 0.11	47.99 $\pm$ 0.16
MFCWMS	<b>44.57 <math>\pm</math> 0.16</b>	<b>73.32 <math>\pm</math> 0.11</b>	<b>47.53 <math>\pm</math> 0.16</b>	<b>48.33 <math>\pm</math> 0.18</b>	<b>76.38 <math>\pm</math> 0.14</b>	<b>51.17 <math>\pm</math> 0.18</b>



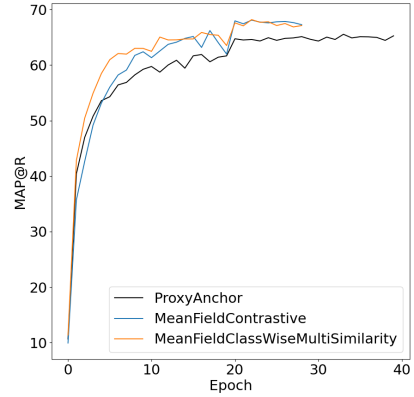
(a) CUB



(b) Cars



(c) SOP



(d) InShop

Figure 1: The test accuracy (MAP@R) plotted against the number of epochs for the (a) CUB, (b) Cars, (c) SOP, and (d) InShop datasets, comparing ProxyAnchor, MFCCont., and MFCWMS.

Table 4: Test accuracies on InShop with the test dataset split into queries and galleries.

Batch size	ProxyAnchor	MFCWMS
30	63.6 $\pm$ 1.4	67.4 $\pm$ 0.2
60	<b>65.7 <math>\pm</math> 0.2</b>	67.6 $\pm$ 0.2
90	65.5 $\pm$ 0.3	67.6 $\pm$ 0.2
120	65.6 $\pm$ 0.3	<b>67.8 <math>\pm</math> 0.2</b>
150	65.5 $\pm$ 0.2	67.0 $\pm$ 0.6
300	64.5 $\pm$ 0.2	67.0 $\pm$ 0.4
500	63.3 $\pm$ 0.2	67.1 $\pm$ 0.1

Table 5: Test accuracies on InShop with the test dataset *not* split into queries and galleries.

Batch size	ProxyAnchor	MFCWMS
30	61.7 $\pm$ 0.6	64.7 $\pm$ 0.1
60	62.8 $\pm$ 0.4	<b>65.1 <math>\pm</math> 0.2</b>
90	<b>62.9 <math>\pm</math> 0.3</b>	65.0 $\pm$ 0.5
120	<b>62.9 <math>\pm</math> 0.3</b>	64.9 $\pm$ 0.5
150	62.7 $\pm$ 0.1	65.0 $\pm$ 0.4
300	61.9 $\pm$ 0.2	64.7 $\pm$ 0.4
500	60.6 $\pm$ 0.2	64.6 $\pm$ 0.1

## References

- [1] Kevin Musgrave, Serge Belongie, and Ser-Nam Lim. A metric learning reality check. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXV 16*, pages 681–699. Springer, 2020.
- [2] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- [3] Sungyeon Kim, Dongwon Kim, Minsu Cho, and Suha Kwak. Proxy anchor loss for deep metric learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3238–3247, 2020.
- [4] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4004–4012, 2016.
- [5] Yair Movshovitz-Attias, Alexander Toshev, Thomas K Leung, Sergey Ioffe, and Saurabh Singh. No fuss distance metric learning using proxies. In *Proceedings of the IEEE international conference on computer vision*, pages 360–368, 2017.
- [6] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1096–1104, 2016.
- [7] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561, 2013.