Multi-layer Radial Basis Function Networks for Out-of-distribution Detection

Supplementary Material

size of 100.

6. Initialization Details

For layer ℓ with N centroids, k-means clustering with N centroids was used to find the initialization positions for the centroids. This can be done quickly for a large batch of data using web-scale k-means clustering with k-means++ initialization [1, 38]. Then for each datapoint in the initialization batch, the distance to its closest centroid is computed, producing a vector in \mathbb{R}^n , where n is the number of datapoints in the batch. Similarly, for each centroid, the distance to its closest datapoint was computed, producing a vector in \mathbb{R}^N . Note that distances were computed using $\|\cdot\|_k$, which is also used in the RBFN forward pass.

To set $\beta_{c(\ell)}^+$, we take inspiration from the Gaussian distribution. In the Gaussian distribution, 95% of datapoints fall within 2 standard deviations of the mean. This can be represented as $2\sigma = d$, or equivalently, $4\sigma^2 = d^2$. Since $\frac{1}{\sigma^2} \propto \beta_{c(\ell)}^+$, we should expect to set $\beta_{c(\ell)}^+$ proportional to $\frac{4}{d^2}$.

Returning to our two vectors in \mathbb{R}^n and \mathbb{R}^N , call the larger of the 95% quantile of each of these metrics d. We can use this d as an estimate of the necessary width of a Gaussian distribution, and thus set $\beta_{c(\ell)}^+ = \frac{4}{d^2}$. Note that we initialize every inverse-width in a layer to the same value, $\frac{4}{d^2}$. In our code, where $\beta_{c(\ell)}^+ = \texttt{softplus}(\beta_{c(\ell)})$, we actually initialize $\beta_{c(\ell)} = \texttt{inverse-softplus}(\frac{4}{d^2})$.

To see this procedure work in code, please see our code in section 9.

7. Moons

To enable development and achieve a simple baseline, we employed a 4-class moons problem. Similar datasets have been used in other papers on uncertainty quantification due to their ease of visualization.

The 4-class moons dataset was generated from scikitlearn's make moons function with noise parameter 0.2, which produces a 2-class moons problem [32]. Two extra classes were generated by shifting the base 2-class moons problem by 2 units in each direction. Overall, the dataset had 1000 training datapoints and 500 testing datapoints split equally among each class. The dataset was normalized prior to training.

We trained a 3-layer MLRBFN with N = 50 centroids in the first and second layers and N = C = 4 centroids in the third layer. A projection dimension of o = 100 was used in the first and second layers. k was set to 2 for all layers. The Adam optimizer was used with a learning rate of 10^{-3} [20]. The network was trained for 250 epochs with a batch

Figure 1 demonstrates the performance of our MLRBFN on the 4-class moons problem. Datapoints are scattered and colored according to their true class. Color in the background represents the class prediction while shading represents the confidence of the network, with white being low

confidence and dark being high confidence. It is clear that the network can classify the dataset correctly (> 99.5% training and testing accuracy) and is only confident near the training data manifold.

8. Results Continued

8.1. Tabular Results

In this section, we include tabular results referred to in the main body of the paper. These can be seen in Table 1, Table 2, Table 3, and Table 4.

In the main body, we perform a deep analysis of Table 1 and Table 2, so here we discuss the other two tables. From Table 3, it is clear that MLRBFNs produce significantly better results than MSP and are competitive with other presented methods. Some other methods, like RMDS, have significantly better performance across all datasets, but these methods can be incorporated with MLRBFNs to produce even better performance.

From Table 4, MLRBFN results are significantly better than MSP: MSP has a 100% false positive rate across all datasets, while MLRBFNs reduce this significantly. Again, MLRBFNs can be combined with other methods presented in the table to produce improved performance.

8.2. OOD Detection with Increasing Number of Classification Layers

Finally, we tested how the number of trainable "head" layers affects OOD detection performance in DNNs and ML-RBFNs. To do this, we trained networks with 1, 2, 3, or 4 head layers on CIFAR10 features extracted by CLIP, and measured AUROC on near-OOD and far-OOD datasets. Results are shown in Figure 5. As can be seen, MLRBFN performance slightly increases as the number of layers increases, while DNNs with ReLU activation functions produce worse performance as the number of layers increases. This may be a sign that feature collapse can occur rapidly, as the DNN's trainable "head" layers map all inputs to one of the training classes even when DNNs use powerful frozen feature extractors. MLRBFNs are able to achieve similar accuracy to DNNs but do not have degrading robustness to OOD inputs.



Figure 5. AUROC for near- and far-OOD tasks on the CIFAR10 dataset. MSP was used for DNN OOD detection. As the number of classification layers increases, MLRBFNs maintain high AUROC for the OOD detection task, while DNNs have decreasing AUROC. This supports MLRBFNs as networks which are naturally capable of detecting OOD inputs.

Table 1. Comparing the performance of MLRBFNs to published benchmarks and existing methods in the OpenOOD v1.5 study. In this table, we use AUROC as the metric for OOD detection. Whenever possible, we report the average number and the corresponding standard deviation obtained from 3 training runs. In almost all cases, MLRBFNs using the ResNet18 (for CIFAR10/CIFAR100/ImageNet200) and ResNet50 (for ImageNet1K) backbones are competitive with the average performance of methods from the OpenOOD v1.5 study. MLRBFNs using DINOv2 and CLIP as feature extracts generally outperform the best performing methods from OpenOOD v1.5.

	CIFAR10			CIFAR100			1	ImageNet200		ImageNet1K		
	Near-OOD	Far-OOD	ID Acc.	Near-OOD	Far-OOD	ID Acc.	Near-OOD	Far-OOD	ID Acc.	Near-OOD	Far-OOD	ID Acc.
- Multi-layer Radial Basis Function Networks												
ResNet	88.88 _(±0.26)	90.01 _(±0.58)	94.63 _(±0.19)	77.30(±0.49)	80.51 _(±1.18)	69.64 _(±0.30)	81.02	89.10	84.50	74.30	86.60	70.72
DINOv2	96.01 _(±0.34)	$97.60_{(\pm 0.64)}$	$97.79_{(\pm 0.04)}$	$89.18_{(\pm 0.73)}$	$93.18_{(\pm 1.17)}$	85.25 _(±0.32)	$90.24_{(\pm 0.22)}$	96.82 _(±0.52)	$94.61_{(\pm 0.37)}$	$78.84_{(\pm 0.25)}$	92.57 _(±0.21)	78.62 _(±0.12)
CLIP	95.75 _(±0.35)	$98.46_{(\pm 0.20)}$	$95.78_{(\pm 0.16)}$	$86.53_{(\pm 0.84)}$	$95.96_{(\pm 1.35)}$	$76.04_{(\pm 0.17)}$	$85.05_{(\pm 0.42)}$	$94.78_{(\pm 0.42)}$	$90.97_{(\pm 0.07)}$	$72.03_{(\pm 0.51)}$	88.50(±0.32)	$72.16_{(\pm 0.11)}$
- OpenOOD v1.5 Be	nchmarks											
Best Post-hoc	90.64 _(±0.20)	93.48 _(±0.24)	95.06 _(±0.30)	81.05 _(±0.07)	82.92(±0.42)	$77.25_{(\pm 0.10)}$	83.69 _(±0.04)	93.90 _(±0.27)	86.37 _(±0.08)	78.17	95.74	76.18
Best Mod. Training	$92.68_{(\pm 0.27)}$	$96.74_{(\pm 0.06)}$	$95.35_{(\pm 0.52)}$	$80.93_{(\pm 0.29)}$	$88.40_{(\pm 0.13)}$	$77.20_{(\pm 0.10)}$	$82.66_{(\pm 0.15)}$	$94.49_{(\pm 0.07)}$	$86.37_{(\pm 0.16)}$	76.52	92.18	76.55
Best Outlier Exp.	94.82 _(±0.21)	$96.00_{(\pm 0.13)}$	$94.95_{(\pm0.04)}$	88.30(±0.10)	$81.41_{(\pm 1.49)}$	$76.84_{(\pm 0.42)}$	$84.84_{(\pm 0.16)}$	$89.02_{(\pm 0.18)}$	$86.12_{(\pm 0.07)}$	N/A	N/A	N/A
Avg. Post-hoc	79.18	83.24	95.06	73.16	76.43	77.25	75.88	84.02	86.37	71.11	85.36	76.18
Avg. Mod. Training	88.06	91.91	94.15	76.28	78.88	73.74	79.44	90.47	85.52	72.96	87.35	75.36
Avg. Outlier Exp.	91.12	93.25	94.12	81.09	78.03	74.84	81.35	87.08	81.44	N/A	N/A	N/A

Table 2. Comparing the performance of MLRBFNs to published benchmarks and existing methods. In this table, we use AUROC as the metric for OOD detection. Whenever applicable, we report the average number and the corresponding standard deviation obtained from 3 training runs. The table demonstrates that MLRBFNs outperform MSP in all cases (emphasized in red), meaning that the outputs of MLRBFNs are better suited for OOD detection than standard DNNs. This result supports our goal of building a neural network architecture which is better able to inherently detect OOD inputs than modern DNNs.

		CIFAR10			CIFAR100			ImageNet200			ImageNet1K		
		Near-OOD	Far-OOD	ID Acc.									
- Multi-lay	er Radial B	asis Function	Networks										
MIDDEN	DINOv2	96.01 _(±0.34)	97.60 _(±0.64)	97.79 _(±0.04)	89.18(±0.73)	93.18 _(±1.17)	85.25 _(±0.32)	90.24 _(±0.22)	96.82 _(±0.52)	94.61 _(±0.37)	78.84 _(±0.25)	92.57 _(±0.21)	78.62 _(±0.12)
WILKDEIN	CLIP	95.75 _(±0.35)	98.46 _(±0.20)	95.78 _(±0.16)	86.53 _(±0.84)	95.96 _(±1.35)	76.04 _(±0.17)	85.05 _(±0.42)	94.78 _(±0.42)	$90.97_{(\pm 0.07)}$	72.03(±0.51)	88.50(±0.32)	72.16(±0.11)
- Foundation	on Model B	enchmarks											
MCD	DINOv2	86.74 _(±0.50)	90.30(±1.59)	98.13 _(±0.09)	83.85 _(±0.46)	87.28 _(±0.88)	87.62(±0.17)	79.09 _(±0.33)	86.25 _(±0.18)	$94.95_{(\pm 0.08)}$	67.80(±0.38)	72.99 _(±0.46)	77.94 _(±0.17)
WISE	CLIP	77.67 _(±0.24)	76.56 _(±1.80)	95.92 _(±0.09)	72.85(±0.67)	72.08(±1.12)	79.14 _(±0.17)	75.55 _(±0.60)	$83.51_{(\pm 0.91)}$	92.73(±0.10)	65.04(±0.32)	69.84 _(±0.67)	$72.21_{(\pm 0.11)}$
ODIN	DINOv2	96.27(±0.27)	$97.16_{(\pm 0.51)}$	$98.13_{(\pm 0.09)}$	85.70 _(±0.54)	89.06(±1.18)	87.62(±0.17)	81.54 _(±0.48)	84.49 _(±0.27)	$94.96_{(\pm 0.08)}$	55.61(±1.70)	49.15(±2.14)	78.42(±0.12)
ODIN	CLIP	83.56(±1.07)	74.32(±2.48)	95.92 _(±0.09)	64.73 _(±1.81)	53.04(±0.79)	79.16(±0.17)	69.18(±1.20)	73.42(±1.88)	$92.74_{(\pm 0.10)}$	61.51 _(±0.94)	63.16(±2.32)	$72.74_{(\pm 0.15)}$
GEN	DINOv2	96.98 _(±0.18)	$97.85_{(\pm 0.30)}$	$98.13_{(\pm 0.09)}$	88.43(±0.20)	$91.73_{(\pm 0.78)}$	87.62 _(±0.17)	88.03(±0.39)	$94.06_{(\pm 0.08)}$	$94.95_{(\pm 0.08)}$	70.17(±0.56)	75.17 _(±0.53)	$77.94_{(\pm 0.17)}$
ULIN	CLIP	89.19 _(±0.23)	$87.01_{(\pm 1.26)}$	$95.92_{(\pm 0.09)}$	74.49 _(±0.65)	$71.62_{(\pm 0.90)}$	$79.14_{(\pm 0.17)}$	$80.68_{(\pm 0.68)}$	$88.12_{(\pm 0.67)}$	$92.73_{(\pm 0.10)}$	68.50 _(±0.40)	$75.33_{(\pm 1.16)}$	$72.21_{(\pm 0.11)}$
SUEID	DINOv2	95.24 _(±0.06)	$97.49_{(\pm 0.06)}$	$98.13_{(\pm 0.09)}$	94.06 _(±0.04)	$95.49_{(\pm 0.39)}$	87.62(±0.17)	91.78 _(±0.04)	$98.30_{(\pm 0.07)}$	$94.95_{(\pm 0.08)}$	81.78 _(±0.17)	$94.83_{(\pm 0.11)}$	77.94 _(±0.17)
SHEII	CLIP	93.53 _(±0.07)	$98.95_{(\pm 0.06)}$	$95.92_{(\pm 0.09)}$	88.74(±0.18)	$97.44_{(\pm 0.20)}$	$79.14_{(\pm 0.17)}$	80.15 _(±0.17)	$89.78_{(\pm 0.48)}$	$92.73_{(\pm 0.10)}$	66.77 _(±0.44)	$79.56_{(\pm 0.73)}$	$72.21_{(\pm 0.11)}$
SHEE	DINOv2	97.62 _(±0.01)	$98.74_{(\pm 0.05)}$	$98.13_{(\pm 0.09)}$	94.89 _(±0.06)	$96.41_{(\pm 0.07)}$	87.62(±0.17)	92.60 _(±0.04)	$97.83_{(\pm 0.08)}$	$94.95_{(\pm 0.08)}$	83.58 _(±0.11)	94.10 _(±0.07)	77.94 _(±0.17)
SHEE	CLIP	94.72 _(±0.05)	$99.10_{(\pm 0.01)}$	$95.92_{(\pm 0.09)}$	89.74(±0.15)	$98.03_{(\pm 0.18)}$	$79.14_{(\pm 0.17)}$	85.57 _(±0.07)	$93.41_{(\pm 0.31)}$	$92.73_{(\pm 0.10)}$	73.97 _(±0.31)	$85.44_{(\pm 0.26)}$	$72.21_{(\pm 0.11)}$
KNN	DINOv2	95.57 _(±0.00)	$96.06_{(\pm 0.00)}$	$98.13_{(\pm 0.09)}$	91.44 _(±0.00)	$86.57_{(\pm 0.00)}$	87.62(±0.17)	$81.51_{(\pm 0.00)}$	$96.20_{(\pm 0.00)}$	$94.95_{(\pm 0.08)}$	77.53 _(±0.00)	$94.65_{(\pm 0.00)}$	77.94 _(±0.17)
KININ	CLIP	90.71 _(±0.00)	$97.82_{(\pm 0.00)}$	$95.92_{(\pm 0.09)}$	82.14(±0.00)	$91.39_{(\pm 0.00)}$	$79.14_{(\pm 0.17)}$	$68.36_{(\pm 0.00)}$	$72.28_{(\pm 0.00)}$	$92.73_{(\pm 0.10)}$	57.80(±0.00)	$61.12_{(\pm 0.00)}$	$72.21_{(\pm 0.11)}$
RMDS	DINOv2	98.27 _(±0.00)	99.19 _(±0.00)	$98.13_{(\pm 0.09)}$	95.16 _(±0.00)	$97.53_{(\pm 0.00)}$	87.62 _(±0.17)	94.77 _(±0.00)	$98.89_{(\pm 0.00)}$	$94.95_{(\pm 0.08)}$	83.79 _(±0.00)	96.31 _(±0.00)	$77.94_{(\pm 0.17)}$
KMDS	CLIP	95.73 _(±0.00)	$97.81_{(\pm 0.00)}$	$95.92_{(\pm 0.09)}$	93.23 _(±0.00)	$98.81_{(\pm 0.00)}$	$79.14_{(\pm 0.17)}$	92.09(±0.00)	$97.88_{(\pm 0.00)}$	$92.73_{(\pm 0.10)}$	78.21 _(±0.00)	$91.50_{(\pm 0.00)}$	$72.21_{(\pm 0.11)}$

Table 3. Comparing the performance of MLRBFNs to published benchmarks and existing methods. In this table, we use AUPR Out as the metric for OOD detection. Whenever applicable, we report the average number and the corresponding standard deviation obtained from 3 training runs.

		CIFAR10			CIFAR100			:	ImageNet200		ImageNet1K		
		Near-OOD	Far-OOD	ID Acc.									
- Multi-lay	er Radial B	asis Function	Networks										
	DINOv2	95.77 _(±0.41)	98.20(±0.60)	97.79 _(±0.04)	86.47 _(±1.13)	93.80 _(±1.23)	85.25 _(±0.32)	93.50(±0.17)	95.83 _(±0.62)	94.61 _(±0.37)	58.82(±1.07)	69.61 _(±0.62)	78.62(±0.12)
MLKBFN	CLIP	95.27 _(±0.51)	98.65 _(±0.30)	95.78 _(±0.16)	83.36 _(±1.41)	95.92 _(±1.70)	76.04 _(±0.17)	88.56(±0.24)	93.55 _(±0.76)	90.97 _(±0.07)	47.53(±0.91)	61.72 _(±1.21)	72.16 _(±0.11)
- Foundatio	on Model B	enchmarks											
MCD	DINOv2	91.17 _(±0.35)	94.80 _(±0.68)	98.13 _(±0.09)	78.91 _(±0.65)	87.86 _(±0.87)	87.62(±0.17)	88.75(±0.23)	88.85 _(±0.09)	94.95 _(±0.08)	43.90(±0.31)	34.34(±0.45)	77.94 _(±0.17)
MSP	CLIP	82.40(±0.15)	85.64 _(±0.97)	95.92 _(±0.09)	65.07 _(±0.48)	75.14 _(±0.76)	79.14 _(±0.17)	83.38(±0.29)	83.36(±0.51)	92.73 _(±0.10)	40.26(±0.24)	29.96(±0.54)	72.21 _(±0.11)
ODIN	DINOv2	95.24 _(±0.48)	96.94 _(±0.68)	98.13(±0.09)	80.47 _(±0.54)	88.82(±1.41)	87.62(±0.17)	84.27(±0.16)	79.72(±0.46)	94.96 _(±0.08)	33.94(±0.89)	16.96(±0.69)	78.42(±0.12)
ODIN	CLIP	80.52 _(±1.62)	77.14(±2.07)	95.92 _(±0.09)	57.75 _(±1.48)	62.89 _(±0.44)	79.16(±0.17)	73.26(±0.85)	66.92 _(±1.71)	92.74 _(±0.10)	36.44 _(±0.97)	23.54(±1.62)	72.74(±0.15)
CEN	DINOv2	96.30(±0.24)	97.93 _(±0.33)	98.13(±0.09)	84.12(±0.32)	91.88 _(±0.94)	87.62(±0.17)	91.45 _(±0.14)	92.34(±0.12)	94.95 _(±0.08)	43.91(±0.81)	33.94(±0.77)	77.94 _(±0.17)
GEN	CLIP	87.69 _(±0.20)	88.66 _(±1.00)	95.92 _(±0.09)	64.77 _(±0.51)	$73.29_{(\pm 0.48)}$	79.14 _(±0.17)	84.16(±0.57)	85.00 _(±1.00)	$92.73_{(\pm 0.10)}$	42.15(±0.45)	35.34(±1.83)	$72.21_{(\pm 0.11)}$
CHEID	DINOv2	95.01 _(±0.08)	98.20(±0.09)	98.13 _(±0.09)	93.60 _(±0.08)	$97.05_{(\pm 0.28)}$	87.62(±0.17)	94.45 _(±0.08)	$97.71_{(\pm 0.18)}$	$94.95_{(\pm 0.08)}$	63.19(±0.55)	$76.19_{(\pm 0.40)}$	$77.94_{(\pm 0.17)}$
SHEIF	CLIP	93.08 _(±0.07)	$98.69_{(\pm 0.05)}$	95.92 _(±0.09)	87.61 _(±0.23)	97.50 _(±0.24)	79.14 _(±0.17)	81.93 _(±0.28)	85.46 _(±0.62)	$92.73_{(\pm 0.10)}$	41.70(±0.37)	37.70 _(±1.12)	$72.21_{(\pm 0.11)}$
SHEE	DINOv2	97.65 _(±0.02)	$99.30_{(\pm 0.05)}$	98.13 _(±0.09)	94.06(±0.09)	97.30 _(±0.06)	87.62(±0.17)	94.44 _(±0.05)	96.62 _(±0.20)	$94.95_{(\pm 0.08)}$	64.35(±0.53)	70.89 _(±0.62)	$77.94_{(\pm 0.17)}$
SHEE	CLIP	94.11 _(±0.06)	$98.76_{(\pm 0.03)}$	95.92 _(±0.09)	88.66 _(±0.24)	$97.64_{(\pm 0.28)}$	79.14 _(±0.17)	86.44 _(±0.26)	90.38 _(±0.46)	$92.73_{(\pm 0.10)}$	47.25 _(±0.39)	46.42 _(±0.73)	$72.21_{(\pm 0.11)}$
KNN	DINOv2	95.37 _(±0.00)	$97.76_{(\pm 0.00)}$	98.13 _(±0.09)	91.23 _(±0.00)	$91.47_{(\pm 0.00)}$	87.62(±0.17)	86.08(±0.00)	$96.14_{(\pm 0.00)}$	$94.95_{(\pm 0.08)}$	60.15(±0.00)	$80.16_{(\pm 0.00)}$	$77.94_{(\pm 0.17)}$
MININ	CLIP	90.51 _(±0.00)	$97.64_{(\pm 0.00)}$	95.92 _(±0.09)	82.59(±0.00)	$91.52_{(\pm 0.00)}$	79.14 _(±0.17)	69.15(±0.00)	$61.81_{(\pm 0.00)}$	$92.73_{(\pm 0.10)}$	35.04(±0.00)	$21.54_{(\pm 0.00)}$	$72.21_{(\pm 0.11)}$
PMDS	DINOv2	97.89 _(±0.00)	99.17 _(±0.00)	98.13 _(±0.09)	94.14 _(±0.00)	98.29 _(±0.00)	87.62(±0.17)	97.16(±0.00)	98.63 _(±0.00)	$94.95_{(\pm 0.08)}$	74.46 _(±0.00)	85.21 _(±0.00)	$77.94_{(\pm 0.17)}$
KMDS	CLIP	94.79 _(±0.00)	$97.08_{(\pm 0.00)}$	$95.92_{(\pm 0.09)}$	92.57 _(±0.00)	$98.91_{(\pm 0.00)}$	$79.14_{(\pm 0.17)}$	94.45 _(±0.00)	$97.27_{(\pm 0.00)}$	$92.73_{(\pm 0.10)}$	55.24 _(±0.00)	$63.57_{(\pm 0.00)}$	$72.21_{(\pm 0.11)}$

Table 4. Comparing the performance of MLRBFNs to published benchmarks and existing methods. In this table, we use FPR@95 (False Positive Rate at 95% True Positive Rate) as the metric for OOD detection. As such, unlike the previous tables, smaller numbers indicate better performance here. Whenever applicable, we report the average number and the corresponding standard deviation obtained from 3 training runs.

		CIFAR10			CIFAR100			ImageNet200			ImageNet1K		
		Near-OOD	Far-OOD	ID Acc.	Near-OOD	Far-OOD	ID Acc.	Near-OOD	Far-OOD	ID Acc.	Near-OOD	Far-OOD	ID Acc.
- Multi-lay	er Radial Basi	is Function Ne	tworks										
MIDDEN	DINOv2	17.99 _(±2.08)	9.51 _(±2.08)	97.79 _(±0.04)	43.39 _(±2.41)	25.31 _(±4.31)	85.25 _(±0.32)	49.55 _(±1.07)	11.71 _(±2.13)	94.61 _(±0.37)	67.08 _(±1.03)	26.92(±1.48)	78.62 _(±0.12)
MLKBFN	CLIP	16.81 _(±1.78)	6.36(±1.30)	95.78 _(±0.16)	46.16(±1.87)	14.21(±3.57)	$76.04_{(\pm 0.17)}$	60.18(±1.89)	22.03(±2.23)	90.97 _(±0.07)	76.68 _(±1.06)	45.61(±2.35)	72.16(±0.11)
- Foundati	on Model Ben	chmarks											
MCD	DINOv2	100.00 _(±0.00)	100.00 _(±0.00)	98.13 _(±0.09)	100.00 _(±0.00)	45.78 _(±1.41)	87.62 _(±0.17)	100.00 _(±0.00)	100.00 _(±0.00)	94.95 _(±0.08)	100.00 _(±0.00)	100.00 _(±0.00)	77.94 _(±0.17)
MSP	CLIP	100.00 _(±0.00)	$100.00_{(\pm 0.00)}$	95.92 _(±0.09)	$100.00_{(\pm 0.00)}$	77.12 _(±7.57)	$79.14_{(\pm 0.17)}$	$100.00_{(\pm 0.00)}$	100.00 _(±0.00)	92.73 _(±0.10)	100.00 _(±0.00)	100.00 _(±0.00)	72.21 _(±0.11)
ODIN	DINOv2	14.11 _(±0.87)	$10.78_{(\pm 1.71)}$	$98.13_{(\pm 0.09)}$	49.59(±2.62)	$33.50_{(\pm 2.68)}$	$87.62_{(\pm 0.17)}$	62.22 _(±2.24)	$47.14_{(\pm 1.25)}$	$94.96_{(\pm 0.08)}$	$87.52_{(\pm 0.61)}$	88.99 _(±1.36)	$78.42_{(\pm 0.12)}$
ODIN	CLIP	55.69 _(±2.25)	$65.67_{(\pm 5.08)}$	$95.92_{(\pm 0.09)}$	86.48(±3.18)	$75.28_{(\pm 1.54)}$	$79.16_{(\pm 0.17)}$	77.12 _(±1.43)	$67.04_{(\pm 2.91)}$	$92.74_{(\pm 0.10)}$	82.37 _(±1.23)	$78.30_{(\pm 1.81)}$	$72.74_{(\pm 0.15)}$
GEN	DINOv2	11.56(±0.89)	8.32(±1.16)	$98.13_{(\pm0.09)}$	42.56(±0.82)	$27.86_{(\pm 2.36)}$	$87.62_{(\pm 0.17)}$	54.84(±3.23)	$24.37_{(\pm 0.28)}$	$94.95_{(\pm 0.08)}$	$76.21_{(\pm 1.17)}$	$67.28_{(\pm 0.99)}$	$77.94_{(\pm 0.17)}$
OLN	CLIP	46.15(±1.00)	$53.18_{(\pm 5.71)}$	$95.92_{(\pm 0.09)}$	67.44 _(±2.49)	$65.24_{(\pm 2.82)}$	$79.14_{(\pm 0.17)}$	66.89 _(±1.49)	$44.34_{(\pm 2.10)}$	$92.73_{(\pm 0.10)}$	$76.47_{(\pm 0.81)}$	$67.22_{(\pm 1.64)}$	$72.21_{(\pm 0.11)}$
SHEIP	DINOv2	18.00(±0.08)	$7.85_{(\pm 0.18)}$	$98.13_{(\pm 0.09)}$	24.63(±0.25)	$14.80_{(\pm 0.71)}$	$87.62_{(\pm 0.17)}$	$43.91_{(\pm 0.28)}$	5.70(±0.18)	$94.95_{(\pm 0.08)}$	$65.63_{(\pm 0.46)}$	$20.64_{(\pm 0.88)}$	$77.94_{(\pm 0.17)}$
SILLI	CLIP	23.36(±0.15)	$4.52_{(\pm 0.20)}$	$95.92_{(\pm 0.09)}$	32.41 _(±0.51)	$10.30_{(\pm 0.83)}$	$79.14_{(\pm 0.17)}$	$60.18_{(\pm 0.39)}$	$32.65_{(\pm 0.81)}$	$92.73_{(\pm 0.10)}$	$80.41_{(\pm 0.52)}$	55.66 _(±1.29)	$72.21_{(\pm 0.11)}$
SHEE	DINOv2	10.09 _(±0.06)	4.22(±0.13)	$98.13_{(\pm 0.09)}$	21.16(±0.04)	11.03 _(±0.27)	87.62 _(±0.17)	40.78 _(±0.13)	$6.19_{(\pm 0.15)}$	$94.95_{(\pm 0.08)}$	$59.54_{(\pm 0.35)}$	$20.29_{(\pm 0.38)}$	$77.94_{(\pm 0.17)}$
STILL	CLIP	18.93 _(±0.13)	3.62(±0.03)	$95.92_{(\pm 0.09)}$	31.17 _(±0.37)	6.79 _(±0.52)	$79.14_{(\pm 0.17)}$	50.88 _(±0.13)	$25.04_{(\pm 0.58)}$	$92.73_{(\pm 0.10)}$	$69.66_{(\pm 0.21)}$	$42.95_{(\pm0.49)}$	$72.21_{(\pm 0.11)}$
WNIN	DINOv2	15.25 _(±0.00)	$10.92_{(\pm 0.00)}$	$98.13_{(\pm 0.09)}$	33.12 _(±0.00)	$27.77_{(\pm 0.00)}$	87.62 _(±0.17)	66.97 _(±0.00)	$19.31_{(\pm0.00)}$	$94.95_{(\pm 0.08)}$	$71.93_{(\pm 0.00)}$	$23.96_{(\pm 0.00)}$	$77.94_{(\pm 0.17)}$
KININ	CLIP-KNN	29.77 _(±0.00)	$8.64_{(\pm 0.00)}$	$95.92_{(\pm 0.09)}$	$44.97_{(\pm 0.00)}$	$23.89_{(\pm 0.00)}$	$79.14_{(\pm 0.17)}$	70.37 _(±0.00)	$47.78_{(\pm 0.00)}$	$92.73_{(\pm 0.10)}$	$83.28_{(\pm 0.00)}$	$63.57_{(\pm 0.00)}$	$72.21_{(\pm 0.11)}$
RMDS	DINOv2	6.41 _(±0.00)	$2.72_{(\pm 0.00)}$	$98.13_{(\pm 0.09)}$	20.59 _(±0.00)	$10.51_{(\pm 0.00)}$	$87.62_{(\pm 0.17)}$	33.93 _(±0.00)	$3.90_{(\pm 0.00)}$	$94.95_{(\pm0.08)}$	$64.98_{(\pm 0.00)}$	$17.07_{(\pm 0.00)}$	$77.94_{(\pm 0.17)}$
KMD8	CLIP	17.33(±0.00)	$8.50_{(\pm 0.00)}$	$95.92_{(\pm 0.09)}$	25.61 _(±0.00)	$4.25_{(\pm 0.00)}$	$79.14_{(\pm 0.17)}$	38.51 _(±0.00)	$8.46_{(\pm 0.00)}$	$92.73_{(\pm 0.10)}$	$66.41_{(\pm 0.00)}$	$30.62_{(\pm 0.00)}$	$72.21_{(\pm 0.11)}$

9. Code

In this section, we provide code for RBF layers of an MLRBFN. Additionally, we provide the code for our binary crossentropy loss function.

```
1
   import torch
2
   import torch.nn as nn
   import torch.nn.functional as F
3
   from torch.nn.modules.lazy import _LazyProtocol
4
   def logsubstractexp(tensor, other):
6
       a = torch.max(tensor, other)
7
       return a + ((tensor - a).exp() - (other - a).exp()).log()
8
   def log_bce_loss(log_y_pred, y, num_classes):
10
       y = F.one_hot(y, num_classes=num_classes).float()
11
       y_flat = y.ravel()
12
       log_y_pred_flat = log_y_pred.ravel() - 1e-6
13
        zero = torch.zeros(log_y_pred_flat.shape).to(y.device)
14
15
16
        not_log_y_pred = logsubstractexp(zero, log_y_pred_flat)
        return torch.mean(-1 * ((y_flat * log_y_pred_flat) + (1 - y_flat) * not_log_y_pred))
17
18
   def webscale_kmeans(data, centroid_shape, k, iter=1):
19
       v = torch.zeros(centroid_shape[0]).to(data.device)
20
        d = torch.zeros(data.shape[0]).int().to(data.device)
21
22
        centroids = torch.empty(centroid_shape).to(data.device)
23
        centroids[0, :] = data[0]
24
       for c in range(1, centroids.shape[0]):
25
           dists = torch.cdist(data, centroids[:c], p=k) ** k
26
           dists = torch.min(dists, -1).values
27
            dists = dists / torch.sum(dists)
28
29
            rand_unif = torch.rand(1).to(data.device)
30
            dists_cumsum = torch.cumsum(dists, 0)
31
            idx = torch.argmax((dists_cumsum > rand_unif).to(torch.long))
32
            centroids[c, :] = data[idx, :] + 1e-5 * torch.randn(data.shape[1]).to(data.device)
33
34
35
       for t in range(iter):
            for i, x in enumerate(data):
36
                x = x[None, :]
37
                dists = torch.linalg.norm(centroids - x, ord=k, dim=1) ** k
38
                idx = torch.argmin(dists)
39
                d[i] = idx
40
41
            for i, x in enumerate(data):
42
                c = d[i]
43
                v[c] += 1
44
                eta = 1 / v[c]
45
46
                centroids[c] = (1 - eta) * centroids[c] + eta * x
47
48
        return centroids
49
   class RBFDepressionLayer(nn.Module):
50
        def __init__(self, input_features, num_centroids, projection_features, k=2, last=False):
51
```

```
super().__init__()
52
53
            self.centroids = nn.Parameter(torch.randn(num_centroids, input_features))
54
            self.init_beta = 1
            self.beta = nn.Parameter(torch.ones(num_centroids,))
55
            self.projections = nn.Parameter(torch.randn(num_centroids, projection_features))
56
            self.k = k
57
58
            self.last = last
59
60
        def forward(self, x, depression, recovery=1.1):
61
            distances = torch.cdist(x, self.centroids, p=self.k) ** self.k
62
            sb = F.softplus(self.beta)
63
            if self.last == False:
65
                scales = (sb / F.softplus(self.init_beta)) * torch.exp(-sb * distances)
66
                depressed_scales = scales * \
67
                     torch.minimum(depression * recovery, torch.tensor(1.0).to(x.device))[:, None]
68
                depress_next = torch.max(depressed_scales, 1).values
69
                x = depressed_scales @ self.projections
70
71
                return x, depress_next, scales
            else:
72
                log_scales = -sb * distances
73
                log_recovery = torch.log(torch.tensor(recovery).to(x.device))
74
                log_depression = torch.log(depression)
75
                x = log\_scales + \setminus
76
                     torch.minimum(log_depression + log_recovery, torch.tensor(0.0).to(x.device))[:, None]
77
                return x, log_depression, log_scales
78
79
    class KMeansRBFDepressionLayer(nn.modules.lazy.LazyModuleMixin, RBFDepressionLayer):
80
        cls_to_become: RBFDepressionLayer
81
        centroids: nn.parameter.UninitializedParameter
82
        beta: nn.parameter.UninitializedParameter
83
84
        def __init__(self, input_features, num_centroids, projection_features, k=2, last=False):
85
            super().__init__(input_features, num_centroids, projection_features, k, last)
86
            self.centroid_shape = self.centroids.shape
87
            self.beta_shape = self.beta.shape
88
89
            self.centroids = nn.parameter.UninitializedParameter()
90
            self.beta = nn.parameter.UninitializedParameter()
91
92
        def initialize_parameters(self, x, depression, recovery=1.1):
93
            self.centroids.materialize(self.centroid_shape)
94
            self.beta.materialize(self.beta_shape)
95
96
97
            x_centroid_loc = x[:x.shape[0]//2, :]
            x_beta = x[x.shape[0]//2:, :]
98
99
100
            cents = webscale_kmeans(x_centroid_loc, self.centroid_shape, self.k, iter=100)
            self.centroids[:] = cents
101
102
            distances = torch.cdist(x_beta, self.centroids, p=self.k) ** self.k
103
104
            distances_dp = torch.quantile(torch.min(distances, 1).values, 0.95)
            distances_cent = torch.quantile(torch.min(distances, 0).values, 0.95)
105
            dist = torch.maximum(distances_dp, distances_cent)
106
107
```

108	<pre>if ((4 / dist)) < 5:</pre>
109	<pre>bs = ((4 / dist)).expml().clamp_min(le-6).clamp_max(le4).log()</pre>
110	else:
111	bs = (4 / dist)
112	self.init_beta = bs
113	<pre>self.beta[:] = bs</pre>