

# Appendix: Detecting Out-of-Distribution Examples with Gram Matrices

Chandramouli S. Sastry<sup>1</sup> Sageev Oore<sup>1</sup>

## A. Schematic Diagram

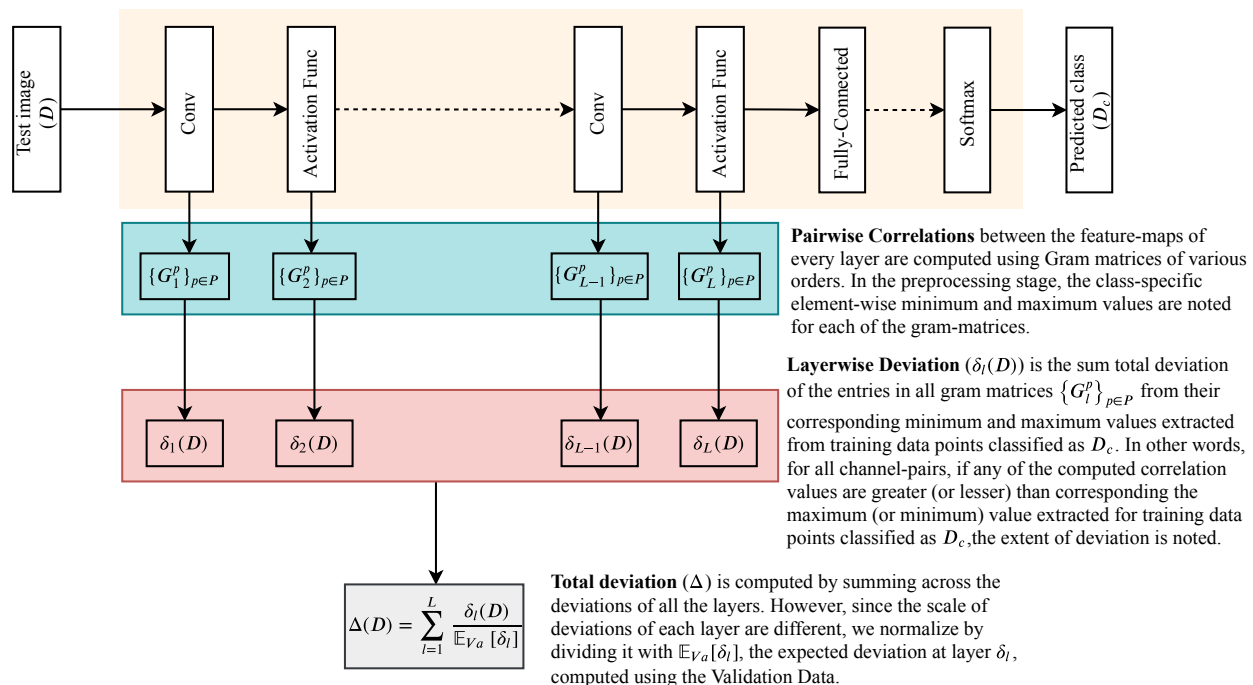


Figure 1: The Schematic Diagram demonstrating the proposed algorithm

## B. Description of OOD Datasets

The following includes the description of the out-of-distribution datasets:

1. **TinyImagenet**, a subset of ImageNet (Russakovsky et al., 2015) images, contains 10,000 test images from 200 different classes. Each image is downsampled to size 32 x 32 and all 10,000 images are used, as given in the opensourced version by (Liang et al., 2018).
2. **LSUN**, the Large-scale Scene UNderstanding dataset (Yu et al., 2015) has 10,000 test images from 10 different scenes. Each image is downsampled to size 32 x 32 and all 10,000 images are used, as given in the opensourced version by (Liang et al., 2018).
3. **iSUN**, a subset of SUN images (Xiao et al., 2010), consists of 8925 images. Each image is downsampled to size 32 x 32 and is used; the downsampled version of the dataset has been opensourced by (Liang et al., 2018).

<sup>1</sup>Dalhousie University/ Vector Institute. Correspondence to: Chandramouli Shama Sastry <chandramouli.sastry@gmail.com>.

4. **SVHN**, the Street View House Numbers dataset (Netzer et al., 2011), involves recognizing digits 0-9 in natural scene images. The test partition consisting of 26,032 images is used.

### C. Detailed Ablation Results

The results in the main paper correspond to the performance obtained when considering:

1. **Feature Set**: all gram matrix entries
2. **Metric**: layerwise deviations computed with respect to the mins and maxs.
3. **Aggregation Scheme**: the total deviation is then computed using Eq 5.

In this section, detailed ablation results are reported by considering other options. Specifically:

1. **Alternate Feature Set**: In addition to considering all gram matrix entries, we consider a proper partition of the gram matrix: strictly diagonal elements, and strictly off-diagonal elements. The diagonal elements correspond to the unary features, while the off-diagonal elements correspond to pairwise features. This can be done by appropriately changing the definition of variable *stat* in Line 7 of Algorithm 1. In these experiments, we consider row-wise sums wherever the size of *stat* is  $O(n^2)$ ; in other words, we consider row-wise sums when considering off-diagonal elements and all gram matrix entries.
2. **Alternate Metric**: An alternative formulation for computing feature-wise deviations can be to compute the deviation from the means using the one-dimensional Mahalanobis distance. In the preprocessing stage, this would be done by storing the *Means* and *Variances* of *stat* (feature-wise) instead of their *Mins* and *Maxs*. Under this new alternative, the function  $\delta$  defined in Eq 3 would be redefined as:

$$\delta(\text{mean, variance, value}) = \frac{(\text{value} - \text{mean})^2}{\text{variance}} \quad (\text{C.1})$$

Accordingly, the layerwise deviation  $\delta_l$  can be defined as:

$$\delta_l(D) = \sum_{p=1}^P \sum_{i=1}^{|\overline{G_l^p(D)}|} \delta \left( \text{Means}[D_c][l][p][i], \text{Variances}[D_c][l][p][i], \overline{G_l^p(D)}[i] \right) \quad (\text{C.2})$$

where  $\overline{G_l^p(D)}$  would correspond to the statistic chosen in the previous step: diagonal entries only, row-wise sums of off-diagonal entries only or row-wise sums of complete gram matrix.

We thus consider 2 options for computing the deviations: the Min/Max method presented in the main paper and the Mean/Variance method (Gaussian) described above. While the Mean/Var assumes each entry of the gram matrix to be normally distributed, the Min/Max assumes that each entry is uniformly distributed between the corresponding extrema and has exponentially decreasing density beyond the extrema:

$$p(\text{value}|\text{min}, \text{max}) = \begin{cases} k & \text{if } \text{min} \leq \text{value} \leq \text{max} \\ k \exp\left(\frac{\text{value}-\text{min}}{|\text{min}|}\right) & \text{if } \text{value} < \text{min} \\ k \exp\left(\frac{\text{max}-\text{value}}{|\text{max}|}\right) & \text{if } \text{value} > \text{max} \end{cases} \quad (\text{C.3})$$

where,  $k = \frac{1}{\text{max}-\text{min} + \frac{1}{|\text{min}|} + \frac{1}{|\text{max}|}}$  makes the above a valid probability density function. Note that both of the proposed density models assume that the entries of the gram matrix are independent of each other. The corresponding  $\delta_l$  can be obtained as the sum of the log probability density estimates.

3. **Alternate Aggregation Scheme**: In order to compute the total deviation  $\Delta$  from the layerwise deviations  $\delta_l$ , we can compute it by following Eq 5 or taking a simple sum as shown:

$$\Delta(D) = \sum_{l=1}^L \delta_l \quad (\text{C.4})$$

We refer to Eq 5 as the normalized estimate and Eq C.4 as the unnormalized estimate.



The proposed Min/Max idea solves problem (a) by employing a weaker metric: deviation from extrema instead of the mean. It can also be said that the Min/Max metric considers a uniform probability density between the extrema. Problem (b), which exists even for this newer metric, is solved by the normalization scheme described in the main paper for computing the sum total deviation.

**Higher Order Gram Matrices** The Min/Max metric is a weak approximation to the true probability density. On conducting a thorough analysis of how the OOD examples were able to fool the metric, it appeared that the intermediate features had several tiny activations that could yield innocuous gram matrix entries. For example, observe in Fig. 2 of the main paper that the Min/Max metric already gets a detection rate close to that of Mahalanobis by using just Order-1 Gram Matrices. Higher-order gram matrices as described in the main paper provide a natural way to mitigate these effects. More importantly, they help in obtaining descriptive summaries of the high-dimensional feature representations through the higher-order non-central moments – of channels and inter-channel hadamard products – contained in them

Notable observations from Figures 2 through 7 (all layers are considered but only one order of gram matrix is considered at a time):

- **Ensemble effect:** In 24/28 cases, higher order gram matrices improve detection rates. Higher order gram matrices help both the Min/Max and the Mean/Var metrics. In most cases, the even powers are more helpful than the odd powers; in some cases, the odd powers are more helpful (Ex: DenseNet/CIFAR-100 vs CIFAR-10). Despite these variations, it is possible to get an ensemble effect by considering all possible powers as demonstrated in the main paper.
- In ResNet: CIFAR-10 vs CIFAR-100 and DenseNet: CIFAR-10 vs CIFAR-100, the higher order gram matrices yield lower detection rates. We find these exceptions interesting, and would like to understand them better in future.

**Summary** The unambiguous message from this ablation study is that the Gram matrix contains useful information which can be used for detecting OOD examples. While the standard Mean/Variance metric does not always work well, the proposed Min/Max metric yields consistent performance competitive with state-of-the-art methods. The use of higher-order Gram matrices further boosts the overall performance. Although the Min/Max method can work very well for "far-from-distribution" examples, it does not work well when a fine grained estimate is needed (for example, CIFAR-10 vs CIFAR-100). We hope the strong empirical proof that Gram matrices contain useful information can motivate the development of OOD detectors with powerful density estimators.

### C.1. Importance of higher-order Gram Matrices

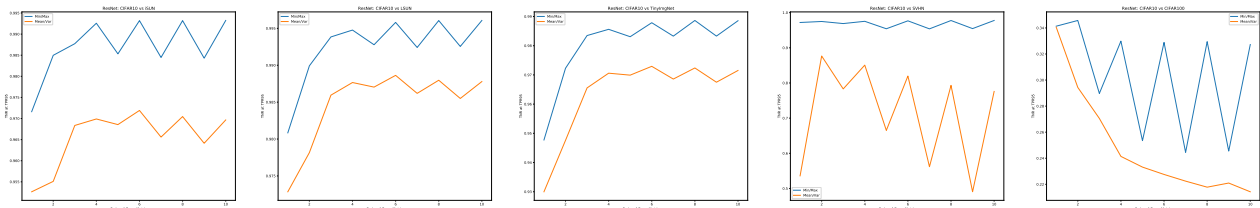


Figure 2: ResNet/CIFAR-10: The TNR at TPR95 trends for Min/Max and Mean/Var as the order of Gram Matrix is varied.

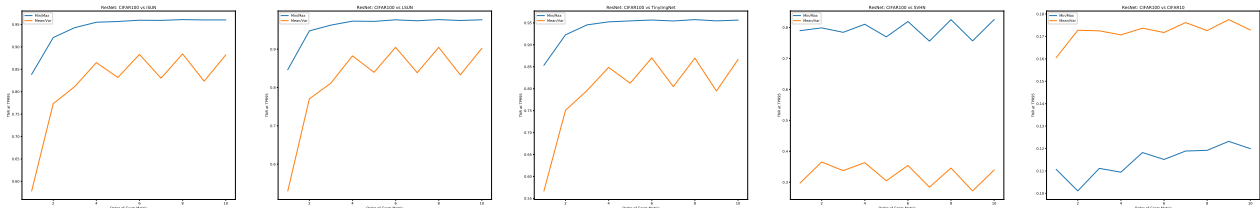


Figure 3: ResNet/CIFAR-100: The TNR at TPR95 trends for Min/Max and Mean/Var as the order of Gram Matrix is varied.

## Appendix: Detecting Out-of-Distribution Examples with Gram Matrices

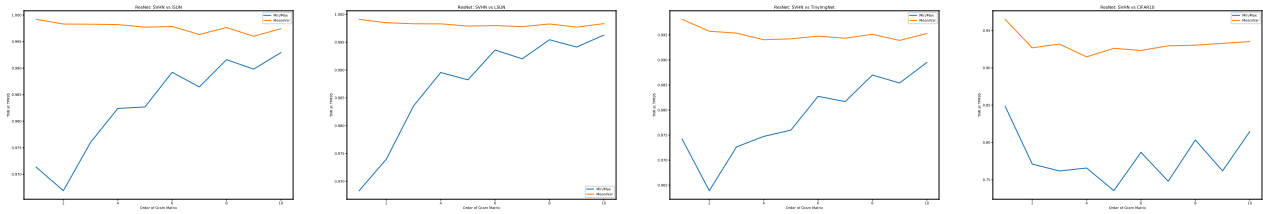


Figure 4: ResNet/SVHN: The TNR at TPR95 trends for Min/Max and Mean/Var as the order of Gram Matrix is varied.

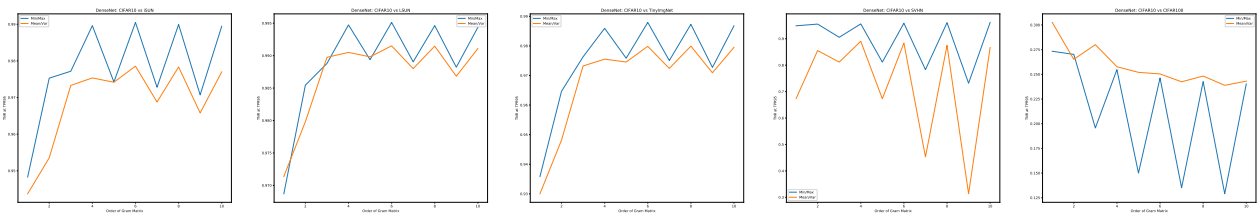


Figure 5: DenseNet/CIFAR-10: The TNR at TPR95 trends for Min/Max and Mean/Var as the order of Gram Matrix is varied.

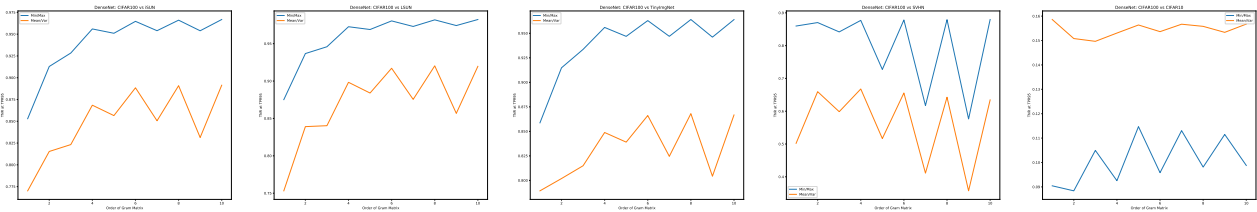


Figure 6: DenseNet/CIFAR-100: The TNR at TPR95 trends for Min/Max and Mean/Var as the order of Gram Matrix is varied.

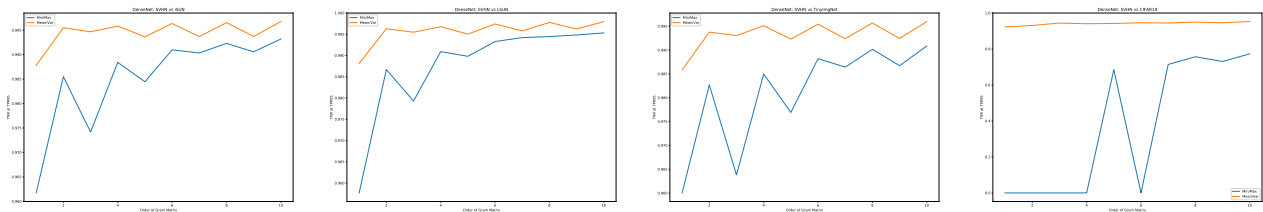


Figure 7: DenseNet/SVHN: The TNR at TPR95 trends for Min/Max and Mean/Var as the order of Gram Matrix is varied.

### C.2. Significance of Depth

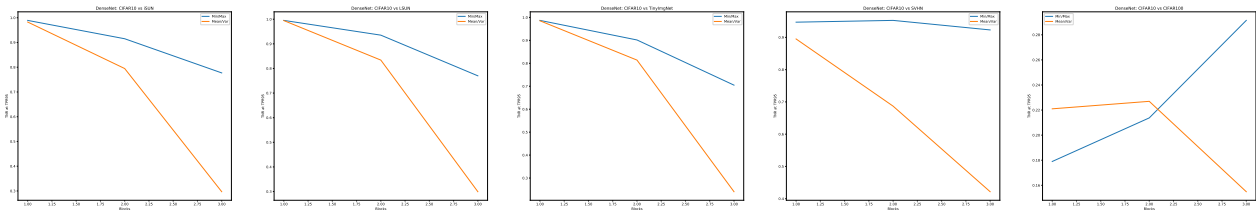


Figure 8: DenseNet/CIFAR-10: The TNR at TPR95 trends for Min/Max and Mean/Var as we go deeper in the network.

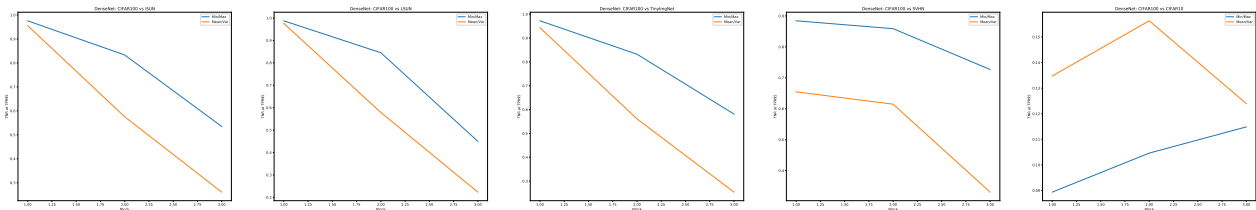


Figure 9: DenseNet/CIFAR-100: The TNR at TPR95 trends for Min/Max and Mean/Var as we go deeper in the network.

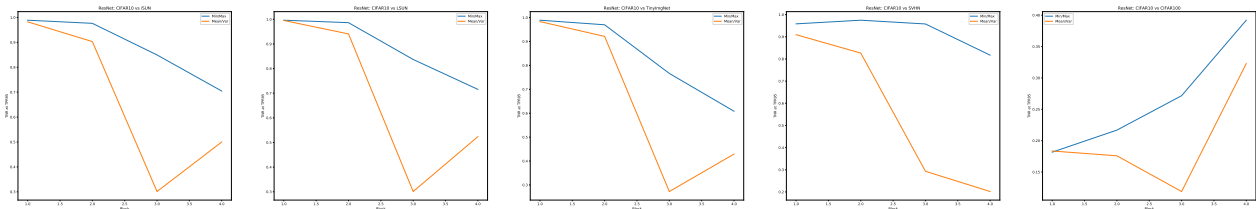


Figure 10: ResNet/CIFAR-10: The TNR at TPR95 trends for Min/Max and Mean/Var as we go deeper in the network.

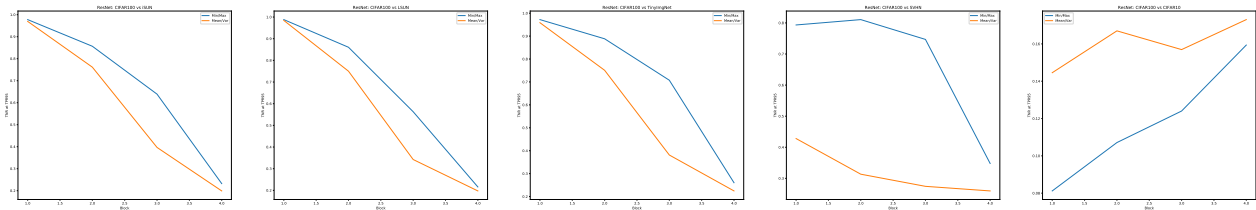


Figure 11: ResNet/CIFAR-100: The TNR at TPR95 trends for Min/Max and Mean/Var as we go deeper in the network.

## Appendix: Detecting Out-of-Distribution Examples with Gram Matrices

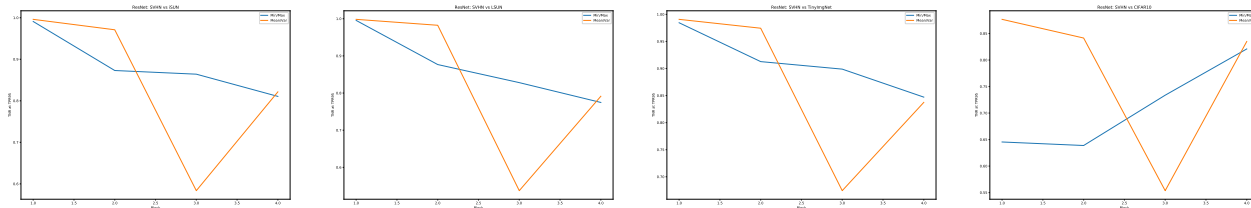


Figure 12: DenseNet/SVHN: The TNR at TPR95 trends for Min/Max and Mean/Var as we go deeper in the network.

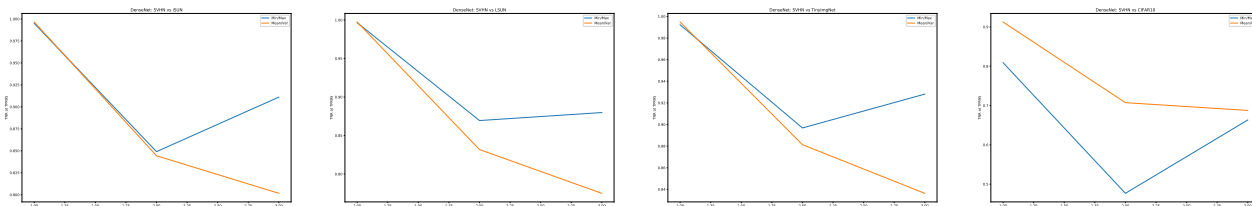


Figure 13: DenseNet/SVHN: The TNR at TPR95 trends for Min/Max and Mean/Var as we go deeper in the network.

### D. Combining OE + Ours

In-dist (WRN 40-2)	OOD	MSP			Ours			Ours + MSP		
		TNR at TPR 95%	AUROC	DTACC	TNR at TPR 95%	AUROC	DTACC	TNR at TPR 95%	AUROC	DTACC
CIFAR-10	iSUN	98.3	99.3	96.9	98.9	99.8	97.8	<b>99.8</b>	99.9	99.0
	LSUN (R)	98.5	99.4	97.0	99.4	99.9	98.4	<b>99.8</b>	99.9	99.1
	LSUN (C)	98.0	99.4	96.9	89.5	97.8	92.5	<b>98.6</b>	99.6	97.3
	TinyImgNet (R)	93.9	98.5	94.6	98.5	99.7	97.6	<b>99.5</b>	99.9	98.5
	TinyImgNet (C)	95.2	98.7	95.2	95.9	99.1	95.7	<b>99.1</b>	99.8	97.8
	SVHN	98.0	99.5	96.9	97.6	99.4	96.8	<b>99.3</b>	99.8	98.2
	CIFAR-100	<b>73.9</b>	94.8	87.9	38.9	80.1	73.3	72.9	93.9	87.0
CIFAR-100	iSUN	50.9	89.8	82.3	<b>96.3</b>	99.1	95.9	95.6	98.9	96.0
	LSUN (R)	58.3	92.0	84.7	<b>98.4</b>	99.6	97.3	97.4	99.3	97.4
	LSUN (C)	69.5	94.0	86.6	69.7	92.6	85.3	<b>83.1</b>	96.3	89.7
	TinyImgNet (R)	36.1	85.1	77.5	<b>96.3</b>	99.1	95.9	92.8	98.2	94.6
	TinyImgNet (C)	41.6	86.3	78.6	<b>90.1</b>	97.7	92.8	87.1	96.9	91.1
	SVHN	56.2	92.5	85.6	84.8	96.5	90.8	<b>85.6</b>	96.8	90.4
	CIFAR-10	<b>17.4</b>	78.4	71.7	7.5	59.3	57.3	16.5	77.7	71.6

Table 2: Table shows results when our method is combined with OE. The experiment was conducted with WideResNet trained with outlier-exposure, open-sourced by (Hendrycks et al., 2019). MSP uses Maximum Softmax Probability; "Ours" refers to the metric  $\Delta$  (Eq. 5); "Ours+MSP" is obtained by using  $\Delta'(x) = \frac{\Delta(x)}{\max_{x \in V_a} \Delta(x)} - \text{MSP}$ .

## E. Few more OOD results

## E.1. Comparing with OE

In-distribution	OOD	OE (Base)	OE	Ours (Base)	Ours
CIFAR-10	Gaussian	85.6	<b>99.3</b>	43.5	<b>100.</b>
	Rademacher	52.4	<b>99.5</b>	48.3	<b>100.</b>
	Blob	83.8	<b>99.4</b>	52.9	<b>99.8</b>
	Texture	57.2	<b>87.8</b>	37.0	85.3
	SVHN	71.2	95.2	45.4	<b>96.1</b>
	LSUN	61.3	87.9	58.2	<b>99.5</b>
CIFAR-100	Gaussian	45.7	87.9	18.2	<b>100.</b>
	Rademacher	61.0	82.9	15.6	<b>100.</b>
	Blob	62.0	87.9	38.4	<b>98.6</b>
	Texture	28.5	45.6	19.9	<b>68.5</b>
	SVHN	30.7	57.1	23.5	<b>85.4</b>
	LSUN	26.0	42.5	18.2	<b>97.2</b>
SVHN	Gaussian	94.6	<b>100.</b>	87.65	<b>100.</b>
	Bernoulli	95.6	<b>100.</b>	92.25	<b>100.</b>
	Blob	96.3	<b>100.</b>	93.35	<b>100.</b>
	Texture	92.8	<b>99.8</b>	72.6	94.9
	Cifar-10	94.0	<b>99.9</b>	73.8	83.0
	LSUN	93.6	<b>99.9</b>	75.7	<b>99.5</b>

Table 3: Comparison of Mean TNR@TPR95 values.

Following (Hendrycks et al., 2019), we created the gaussian, rademacher, blob and bernoulli synthetic datasets. Their descriptions are as follows: *Gaussian* anomalies have each dimension i.i.d. sampled from an isotropic Gaussian distribution. *Rademacher* anomalies are images where each dimension is -1 or 1 with equal probability, so each dimension is sampled from a symmetric Rademacher distribution. *Bernoulli* images have each pixel sampled from a Bernoulli distribution if the input range is [0, 1]. *Blobs* data consist of algorithmically generated amorphous shapes with definite edges. *Textures* is a dataset of describable textural images (Cimpoi et al., 2014).

## E.2. Comparing with DPN, VD and Semantic.

OOD	Method	TNR @ TPR95	AUROC	Detection Accuracy
LSUN	DPN	42.60	90.20	79.50
	VD	92.30	98.30	94.10
	Baseline	49.80	91.00	85.30
	ODIN	82.10	94.10	86.70
	Mahalanobis	98.80	99.70	97.70
	<b>Ours</b>	<b>99.85</b>	<b>99.89</b>	<b>98.66</b>
Tiny ImgNet	DPN	71.60	93.00	86.40
	VD	82.90	96.80	91.30
	Baseline	41.00	91.00	85.10
	ODIN	67.90	94.00	86.50
	Mahalanobis	97.10	99.50	96.30
	<b>Ours</b>	<b>99.48</b>	<b>99.72</b>	<b>97.82</b>
SVHN	DPN	79.90	95.90	87.30
	VD	71.30	93.20	86.40
	Baseline	50.50	89.90	85.10
	ODIN	70.30	96.70	91.10
	Mahalanobis	87.80	99.10	95.80
	<b>Ours</b>	<b>98.14</b>	<b>99.50</b>	<b>96.71</b>

(a) ResNet/CIFAR-10

OOD	Method	TNR @ TPR95	AUROC	Detection Accuracy
iSUN	Semantic	41.60	85.20	88.40
	VD	80.20	94.20	87.80
	Baseline	16.89	75.80	70.11
	ODIN	45.21	85.48	78.47
	Mahalanobis	89.91	97.91	93.05
	<b>Ours</b>	<b>95.12</b>	<b>98.9</b>	<b>95.18</b>
LSUN	Semantic	20.50	79.00	57.80
	VD	85.50	95.90	90.40
	Baseline	18.80	75.80	69.90
	ODIN	23.20	85.60	78.30
	Mahalanobis	90.89	98.2	93.5
	<b>Ours</b>	<b>97.14</b>	<b>99.28</b>	<b>96.19</b>
Tiny ImgNet	Semantic	37.60	83.10	75.60
	VD	83.70	95.30	89.70
	Baseline	20.40	77.20	70.80
	ODIN	36.1	87.6	80.1
	Mahalanobis	90.92	98.20	93.30
	<b>Ours</b>	<b>95.12</b>	<b>98.97</b>	<b>95.13</b>

(b) ResNet/CIFAR-100

Table 4: We compare our method with DPN, VD and Semantic by reporting results where available.



### E.3. Results for Fully-connected Networks

Architecture	OOD	Method	TNR @ TPR95	AUROC	Detection Accuracy
300	KMNIST	Baseline	47.66	73.96	73.91
		<b>Ours</b>	<b>98.57</b>	<b>99.66</b>	<b>97.37</b>
	Fashion-MNIST	Baseline	44.93	66.93	71.07
		<b>Ours</b>	<b>93.51</b>	<b>98.64</b>	<b>94.36</b>
300-150	KMNIST	Baseline	59.79	75.17	79.49
		<b>Ours</b>	<b>97.8</b>	<b>99.4</b>	<b>96.55</b>
	Fashion-MNIST	Baseline	70.73	77.10	83.00
		<b>Ours</b>	<b>95.2</b>	<b>99.00</b>	<b>95.17</b>
300-150-50	KMNIST	Baseline	70.4	79.75	83.38
		<b>Ours</b>	<b>97.5</b>	<b>99.11</b>	<b>96.4</b>
	Fashion-MNIST	Baseline	73.92	76.54	84.67
		<b>Ours</b>	<b>95.7</b>	<b>98.94</b>	<b>95.48</b>

Table 5: The method even works quite well with a fully-connected neural network trained on MNIST. The results are shown for 300-unit single layer MLP, 300-150 two-layer MLP and 300-150-50 MLP.

### F. SVHN images



Figure 14: Some images selected from the test partition of SVHN which have unusual feature correlations as determined by our method. Some images we found interesting include what appears to be a porch lamp (Row 2 Col 5) and an 8 inside a 0 (Row 2 Col 3).

### References

- M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, pages 3606–3613, 2014. doi: 10.1109/CVPR.2014.461. URL <https://doi.org/10.1109/CVPR.2014.461>.
- D. Hendrycks, M. Mazeika, and T. G. Dietterich. Deep anomaly detection with outlier exposure. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019, 2019*. URL <https://openreview.net/forum?id=HyxCxhRcY7>.
- S. Liang, Y. Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, 2018*. URL <https://openreview.net/forum?id=H1VGkIxRZ>.

- Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011. URL [http://ufldl.stanford.edu/housenumbers/nips2011\\_housenumbers.pdf](http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf).
- O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, Dec 2015. ISSN 1573-1405. doi: 10.1007/s11263-015-0816-y. URL <https://doi.org/10.1007/s11263-015-0816-y>.
- J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. SUN database: Large-scale scene recognition from abbey to zoo. In *The Twenty-Third IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2010, San Francisco, CA, USA, 13-18 June 2010*, pages 3485–3492, 2010. doi: 10.1109/CVPR.2010.5539970. URL <https://doi.org/10.1109/CVPR.2010.5539970>.
- F. Yu, Y. Zhang, S. Song, A. Seff, and J. Xiao. LSUN: construction of a large-scale image dataset using deep learning with humans in the loop. *CoRR*, abs/1506.03365, 2015. URL <http://arxiv.org/abs/1506.03365>.