
Improving the out-of-distribution performance of score-based generative models via self-supervision

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this work, we first examine the efficacy of score-based generative models (SGMs)
2 for out-of-distribution (OOD) detection. We show previously proposed OOD de-
3 tection metrics based on SGMs fail to address OODs that share similar textures but
4 different object shapes. Based on the observation, we construct RotNCSN, a novel
5 OOD detection method based-on the score matching and data augmentation. Rot-
6 NCSN first applies random rotation to the perturbed data and forces its output to be
7 rotation-invariant. Therefore, RotNCSN becomes more shape-aware. Experiment
8 results show that RotNCSN consistently improves over the baseline metric based
9 on the SGMs. Furthermore, RotNCSN also achieves competitive OOD detection
10 performance in the FashionMNIST domain.

11 1 Introduction

12 Score-based generative models (SGMs) Song and Ermon [2019], Ho et al. [2020], Song et al. [2021]
13 have emerged as a promising method for deep generative modelling on various domains Dhariwal and
14 Nichol [2021], Xu et al. [2022] due to its competitive performance and stable training. Furthermore,
15 they have been successfully applied in various image-based subtasks, including stroke-based editing
16 Meng et al. [2022], super-resolution Hoogeboom et al. [2022], and segmentation Baranchuk et al.
17 [2022]. However, most of the applications are based on image generation and relatively little work has
18 been devoted to applying SGMs to hypothesis testing, including out-of-distribution (OOD) detection.
19 OOD detection aims to design a reliable metric that discriminates the given distribution from the
20 others. Since generative models naturally model in-distribution images, they are widely applied for
21 OOD detection Havtorn et al. [2021], Xiao et al. [2020]. However, there are few investigations on the
22 potential of SGM in OOD detection.

23 In this paper, we first question whether previous OOD detection metrics based on SGMs determine
24 data based on its object shape. For example, in the work of Yang et al. [2021], OOD detection
25 methods based on the classifier trained in the CIFAR-10 cat and dog images assign higher confidence
26 to the CIFAR-10 data that come from different classes than the dog image from the ImageNet data
27 with negligible covariate shift. We take this analog to an unsupervised setting and design OOD data
28 that share a similar texture to the in-distribution data. OOD detection metrics based on the norm of
29 the score function Mahmood et al. [2021] and reconstruction loss are vulnerable to such OOD data
30 that share a similar texture.

31 To overcome such issues, we propose RotNCSN, an OOD detection method that integrates score
32 matching with data augmentation. RotNCSN is trained to be rotation-invariant and therefore be more
33 shape-aware. We test RotNCSN in various OOD detection benchmarks. Compared to the previous
34 score matching-based methods, RotNCSN achieves considerable improvement. Moreover, RotNCSN
35 shows competitive performance against state-of-the-art baselines on unsupervised OOD detection.

36 **2 Background**

37 In this section, we introduce the formulation of OOD detection. Furthermore, we study examples of
 38 score-based methods applied to OOD detection.

39 **2.1 OOD detection**

40 In OOD detection, we want to distinguish the given distribution \mathcal{D} from the others. Therefore, given
 41 $\mathcal{D}_{\text{train}} \subset \mathcal{D}$, we design a metric f that can discriminate samples from $\mathcal{D}_{\text{test}} \subset \mathcal{D}$ from the other outlier
 42 distributions. We use binary hypothesis testing measures to evaluate the model; area under the ROC
 43 curve (AUROC) or detection accuracy.

44 In unsupervised OOD detection where only the data information is applied, generative models have
 45 been widely used to extract reliable metrics. In the beginning, the likelihood or discriminator output
 46 of the generative models, including VAE Kingma and Welling [2014], GAN Goodfellow et al. [2014],
 47 and GLOW Kingma and Dhariwal [2018], is used. However, Nalisnick et al. [2019] found that such
 48 metrics can be vulnerable to distinguishing simple OODs, such as SVHN from CIFAR-10 Krizhevsky
 49 and Hinton [2009]. Various methods are proposed to explain such phenomena; complexity of the
 50 data Choi and Chung [2020], Serrà et al. [2020], overfitting to low-level features Havtorn et al.
 51 [2021], Kirichenko et al. [2020], Schirrmeister et al. [2020], Zhang et al. [2021], and overfitting into
 52 backgrounds ren et al. [2019]. While SGMs do considerably better in such dataset Mahmood et al.
 53 [2021], we show they are still vulnerable to OODs generated by geometrical transformations.

54 **2.2 SGMs and application to OOD detection**

55 In this section, we introduce MSMA Mahmood et al. [2021], an OOD detection method that utilizes
 56 the score function of SGM. MSMA use NCSN Song and Ermon [2019] as a base model for training
 57 the score function. NCSN utilizes a score function $s_\theta(x, l)$ that takes perturbed data and the degree
 58 of perturbation as input. NCSN trains to match the gradient of log-likelihood of the perturbed data in
 59 multiple scales $(\sigma_i)_{i=1}^L$. We follow the denoising score matching version of NCSN that is optimized
 60 to minimize the loss function below.

$$\frac{1}{2L} \sum_{i=1}^L \sigma_i^2 \mathbb{E}_{x \in \mathcal{D}_{\text{train}}} \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma_i I)} \left[\left\| s_\theta(\tilde{x}, \sigma_i) + \frac{\tilde{x} - x}{\sigma_i^2} \right\|_2^2 \right] \quad (1)$$

61 MSMA utilizes the normality vector $f_{\text{msma}}(x)$ for the observed data x , which is based on the multi-
 62 scale score function. f_{msma} is a L -dimensional vector where each indice is defined as below.

$$f_{\text{msma}}(x)_i = \|\sigma_i s_\theta(x, \sigma_i)\| \quad (2)$$

63 MSMA trains an unsupervised one-class classification model based on the normality vector and
 64 uses its likelihood to detect OOD data. Alongside the score function, we also explore the choice of
 65 normality vector f_{rec} based on the reconstruction error of the perturbation.

$$f_{\text{rec}}(x)_i = \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma_i I)} \left\| s_\theta(\tilde{x}, \sigma_i) + \frac{\tilde{x} - x}{\sigma_i^2} \right\| \quad (3)$$

66 While computing the expectation may require multiple iterations, we found that the number of
 67 iterations does not affect the normality vector much. MSMA learns an unsupervised model (e.g
 68 GMM) over the normality vector of training data to extract scalar metrics. While MSMA shows
 69 state-of-the-art performance in some OOD detection tasks Mahmood et al. [2021], we show that they
 70 show underwhelming performance in detecting various OODs.

Method	CIFAR-100	rot90	rot180	rot270	Patch shuffle
f_{msma}	0.615	0.554	0.542	0.562	0.737
f_{rec}	0.607	0.556	0.546	0.570	0.737

Table 1: AUROC of previously proposed SGM-based OOD detection methods trained in the CIFAR-10 dataset tested in the proposed OOD dataset. All the methods show underwhelming performance on the OOD data.

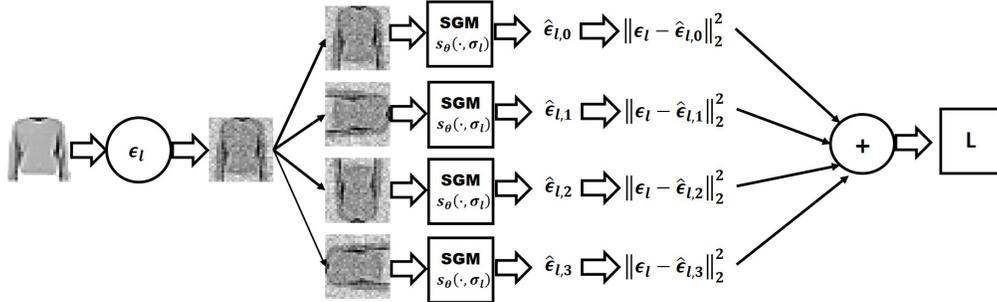


Figure 1: Visual schematic of RotNCSN. RotNCSN forces the output to be rotation-invariant concerning the perturbed data.

71 3 Methodology

72 3.1 Proposed Method

73 3.2 Motivation

74 We first ask our motivating question: do SGMs recognize in-distribution data via object shape instead
 75 of texture? For example, if a classifier model changes its prediction when the texture changes, the
 76 model is likely to predict the data by the texture, not by the object shape Geirhos et al. [2019]. Since
 77 we are dealing with the unsupervised setting, we test SGMs against the model that shares similar
 78 backgrounds but differs in shape. Specifically, we test the MSMA Mahmood et al. [2021] method
 79 on the NCSN trained in the CIFAR-10 Krizhevsky and Hinton [2009] dataset against the following
 80 OODs.

81 • **CIFAR-100** is known as near-OOD data for CIFAR-10 since they share similar textures. Since
 82 CIFAR-100 and CIFAR-10 are both subsets of the 80-million image dataset, an OOD detector should
 83 be aware of class-wise discriminability. This is challenging for an unsupervised OOD detection
 84 setting.

85 • **Rotation** is also a plausible OOD to check the dependence of metric to object orientation. In the
 86 CIFAR-10 dataset, most data share a similar shape orientation. For example, there is no deer standing
 87 upside-down in the CIFAR-10 dataset. Therefore, we regard rotated data as OOD data Gossweiler
 88 et al. [2009]. We test the CIFAR-10 dataset that rotated 90, 180, and 270 degrees counter-clockwise
 89 and refer to them as rot90, rot180, and rot270.

90 • **Patch shuffling** extracts the patch from the image and shuffles the order of the patch to construct
 91 the OOD data. We divide the image into 16×8 -sized patches and shuffle them in random order.
 92 This operation relatively destroys the object’s shape compared to the texture Noroozi and Favaro
 93 [2016].

94 We now report the performance of GMMs trained on f_{msma} Mahmood et al. [2021] and f_{rec} in the
 95 following OODs in Table 1. Both metrics show underwhelming performance in detecting OOD
 96 images that share similar textures. We further show that the trend is consistent when evaluated in the
 97 alternative dataset, SVHN. We refer the results to Appendix.

98 We now introduce our proposed scheme RotNCSN. In the training phase, RotNCSN applies random
 99 rotation $r \in \mathcal{R} = \{\text{Rot}(X, 0), \text{Rot}(X, 90), \text{Rot}(X, 180), \text{Rot}(X, 270)\}$. Then, RotNCSN minimizes
 100 the loss function below.

OOD	RotNCSN (ours)	f_{msma}	f_{rec}	LR	IC	Likelihood Ratio
EMNIST	0.982	0.961	0.937			
MNIST	0.994	0.828	0.842	0.967	0.946	0.924
NotMNIST	0.978	0.932	0.892	1.0	0.923	0.996
KMNIST	0.988	0.901	0.893	0.983	0.708	0.983

OOD	RotNCSN	f_{msma}	f_{rec}
CIFAR-100	0.678	0.615	0.607

Table 2: AUROC of RotNCSN and other OOD detection metrics on FashionMNIST (**up**) and CIFAR-10 (**down**) in-distribution datasets. Results on LR, IC, and Likelihood Ratio are from Xiao et al. [2020]

$$\frac{1}{2L} \sum_{i=1}^L \sigma_i^2 \mathbb{E}_{x \in D_{\text{train}}} \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma_i I)} \mathbb{E}_{r \in \mathcal{R}} \left[\left\| s_{\theta} (r(\tilde{x}), \sigma_i) + \frac{\tilde{x} - x}{\sigma_i^2} \right\|_2^2 \right] \quad (4)$$

101 The score function of RotNCSN outputs the original noise matrix instead of the rotated noise matrix.
102 Therefore, RotNCSN naturally learns to discriminate in-distribution data from the rotated data.
103 Moreover, we expect RotNCSN to be more shape-aware since rotation-based self-supervised learning
104 methods show efficacy in various downstream tasks Gidaris et al. [2018], Hendrycks et al. [2019].

105 When new sample x is given, RotNCSN outputs the normality vector $f_{\text{rot}}(x)$ as follows.

$$f_{\text{rot}}(x)_i = \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma_i I)} \sum_{r \in \mathcal{R}} \left\| s_{\theta} (r(\tilde{x}), \sigma_i) + \frac{\tilde{x} - x}{\sigma_i^2} \right\|_2 \quad (5)$$

106 RotNCSN then trains a unsupervised model (e.g. GMM) over the extracted normality vector from the
107 training data. This is consistent to the evaluation of Mahmood et al. [2021].

108 4 Discussion

109 We first evaluate the performance of RotNCSN in the following OOD detection tasks. We refer
110 training of the RotNCSN in the Appendix.

111 **FashionMNIST**: we evaluate RotNCSN trained in FashionMNIST Xiao et al. [2017] dataset against
112 the various OOD dataset; EMNIST, MNIST LeCun et al. [2010], NotMNIST, and KMNIST Clanuwat
113 et al. [2018].

114 **CIFAR-10**: we evaluate RotNCSN trained in the CIFAR-10 dataset against the challenging CIFAR-
115 100 dataset. Since the two datasets show similar textures, we expect the conventional normality vector
116 to struggle in this task.

117 We sample \tilde{x} once w.r.t each x while we compute $f_{\text{rot}}(x)$. For the baseline, we compare with
118 conventional normality vector f_{msma} and f_{recon} . We directly use the trained score model of Mahmood
119 et al. [2021] for evaluating the baseline method. Since all score-based methods output the normality
120 vector, we train GMM over the normality vector of the train data to output the explicit likelihood.
121 Furthermore, as a competitive unsupervised OOD detection baseline, we compare the results of Xiao
122 et al. [2020], Serrà et al. [2020], and ren et al. [2019] and refer to them as LR, IC, and Likelihood
123 Ratio. For the result of LR and Likelihood Ratio, we refer to the result of Xiao et al. [2020] that use
124 VAE.

125 We show the OOD detection result in Table 2. In the FashionMNIST domain, our proposed RotNCSN
126 consistently improves over the conventional NCSN-based OOD detection metrics (MSMA, Recon).
127 Furthermore, RotNCSN is competitive with various OOD detection methods. Finally, in the CIFAR-
128 10 domain, RotNCSN improves over MSMA in the challenging CIFAR-100 detection task. We leave
129 further analysis in the Appendix.

130 **References**

- 131 D. Baranchuk, A. Voynov, I. Rubachev, V. Khrulkov, and A. Babenko. Label-efficient semantic
132 segmentation with diffusion models. In *ICLR*, 2022.
- 133 S. Choi and S.-Y. Chung. Novelty detection via blurring. In *ICLR*, 2020.
- 134 T. Clanuwat, M. Bober-Irizar, A. Kitamoto, A. Lamb, K. Yamamoto, and D. Ha. Deep learning for
135 classical japanese literature. *arXiv preprint arXiv:1812.01718*, 2018.
- 136 G. Cohen, S. Afshar, J. Tapson, and A. van Schaik. Emnist: an extension of mnist to handwritten
137 letters. *arXiv preprint arXiv:1702.05373*, 2017.
- 138 P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*, 2021.
- 139 R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, and W. Brendel. Imagenet-trained
140 cnns are biased toward texture; increasing shape bias improves accuracy and robustness. In *ICLR*,
141 2019.
- 142 S. Gidaris, P. Singh, and N. Komodakis. Unsupervised representation learning by predicting image
143 rotations. In *ICLR*, 2018.
- 144 I. J. Goodfellow, J. Pouget-Abadie, B. X. Mehdi Mizra, D. Warde-Farley, S. Ozair, A. Courville, and
145 Y. Bengio. Generative adversarial networks. In *NeurIPS*, 2014.
- 146 R. Gossweiler, M. Kamvar, and S. Baluja. What’s up captcha: a captcha based on image orientation.
147 In *WWW*, 2009.
- 148 J. D. Havtorn, J. Frellsen, S. Hauberg, and L. Maaløe. Hierarchical vaes know what they don’t know.
149 In *ICML*, 2021.
- 150 D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song. Using self-supervised learning can improve
151 model robustness and uncertainty. In *NeurIPS*, 2019.
- 152 J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- 153 E. Hooeboom, A. A. Gritsenko, J. Bastings, B. Poole, R. van den Berg, and T. Salimans. Autore-
154 gressive diffusion models. In *ICLR*, 2022.
- 155 D. P. Kingma and P. Dhariwal. Glow: generative flow with invertible 1x1 convolutions. In *NeurIPS*,
156 2018.
- 157 D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014.
- 158 P. Kirichenko, P. Izmailov, and A. G. Wilson. Why normalizing flows fail to detect out-of-distribution
159 data. In *NeurIPS*, 2020.
- 160 A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Technical*
161 *Report*, 2009.
- 162 Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. 2010.
- 163 A. Mahmood, J. Oliva, and M. Styner. Multiscale score matching for out-of-distribution detection. In
164 *ICLR*, 2021.
- 165 C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon. Sdedit: guided image synthesis
166 and editing with stochastic differential equations. In *ICLR*, 2022.
- 167 E. Nalisnick, A. Matsukawa, Y. W. Teh, D. Gorur, and B. Lakshminarayanan. Do deep generative
168 models know what they don’t know? In *ICLR*, 2019.
- 169 M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles.
170 In *ECCV*, 2016.
- 171 J. ren, P. J. Liu, E. Fertig, J. Snoek, R. Poplin, M. A. DePristo, J. V. Dillon, and B. Lakshminarayanan.
172 Likelihood ratios for out-of-distribution detection. In *NeurIPS*, 2019.

- 173 R. T. Schirrmeister, Y. Zhou, T. Ball, and D. Zhang. Understanding anomaly detection with deep
174 invertible networks through hierarchies of distributions and features. In *NeurIPS*, 2020.
- 175 J. Serrà, D. Álvarez, V. Gómez, O. Slizovskaia, J. F. Núñez, and J. Luque. Input complexity and
176 out-of-distribution detection with likelihood-based generative models. In *ICLR*, 2020.
- 177 Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution. In
178 *NeurIPS*, 2019.
- 179 Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative
180 modeling through stochastic differential equations. In *ICLR*, 2021.
- 181 H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dastaset for benchmarking machine
182 learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 183 Z. Xiao, Q. Yan, and Y. Amit. Likelihood regret: an out-of-distribution detection for variational
184 autoencoder. In *NeurIPS*, 2020.
- 185 M. Xu, L. Yu, Y. Song, C. Shi, S. Ermon, and J. Tang. Geodiff: a geometric diffusion model for
186 molecular conformation generation. In *ICLR*, 2022.
- 187 J. Yang, H. Wang, L. Feng, X. Yan, H. Zheng, W. Zhang, and Z. Liu. Semantically coherent
188 out-of-distribution detection. In *ICCV*, 2021.
- 189 M. Zhang, A. Zhang, and S. McDonagh. On the out-of-distribution generalization of probabilistic
190 image modelling. In *NeurIPS*, 2021.



Figure 2: Examples of the tested OODs. (a): Original CIFAR-10 image. (b),(c),(d): rotated CIFAR-10 images. (e): patch-shuffled CIFAR-10 image. (f): CIFAR-100 images.

Method	rot90	rot180	rot270	Patch shuffle
f_{msma}	0.635	0.520	0.630	0.897
f_{rec}	0.640	0.527	0.634	0.865

Table 3: AUROC of previously proposed SGM-based OOD detection methods trained in the SVHN dataset tested in the proposed OOD dataset.

191 A Appendix

192 A.1 Further results on the motivation

193 We further provide the visualization and additional results that support our hypothesis. First, we
 194 provide a visualization of the OOD data used in Figure 2. We further test our motivation in the
 195 alternative SVHN dataset. We refer the result to Table 3. Similar to CIFAR-10, SGM trained in the
 196 SVHN dataset also struggles to detect rotated OODs although they are not in the training dataset.
 197 Nevertheless, SGM trained in the SVHN dataset is more robust to patch shuffling compared to the
 198 CIFAR-10 dataset.

199 A.2 Experiment settings

200 In training the RotNCSN, we do not change any training details (e.g noise scale) except the training
 201 epoch. We train RotNCSN for 400000 steps in the FashionMNIST dataset and 600000 steps in the
 202 CIFAR-10 dataset.

203 A.3 Further analysis

204 In this section, we further analyze the performance of RotNCSN compared to the NCSN. We visualize
 205 top-9 samples that RotNCSN and NCSN output the lowest likelihood on the FashionMNIST domain
 206 in Figure 3. While NCSN assigns higher uncertainty to relatively complex data, RotNCSN assigns
 207 higher uncertainty to data with a possible anomaly. For example, in the third row of the left of Figure
 208 3, we observe cracks in the object.

209 We further analyze the effect of each noise level for OOD detection. While our method originates on
 210 the multi-scale vector, we extract the score of each noise level and use them independently for OOD
 211 detection. We simply use the distance from the mean of the training dataset’s score as the detection
 212 metric. Instead of AUROC, we plot TNR at 95% TPR for the distinct visualization.

213 We report the result in Figure 4. Since we did not train over GMM nor use multi-scale score matching,
 214 the performance is less than the reported one in Table 2. The OOD detection performance increases
 215 while the noise level increases generally. However, the performance diminishes after 0.359. One
 216 hypothesis to explain this behavior is that reconstructing the original image at a high noise level may
 217 be an ill-posed problem and therefore become unsuitable for OOD detection.



Figure 3: Top-9 samples from the Fashion-MNIST dataset where RotNCSN (**left**) and NCSN (**right**) assign highest uncertainty.

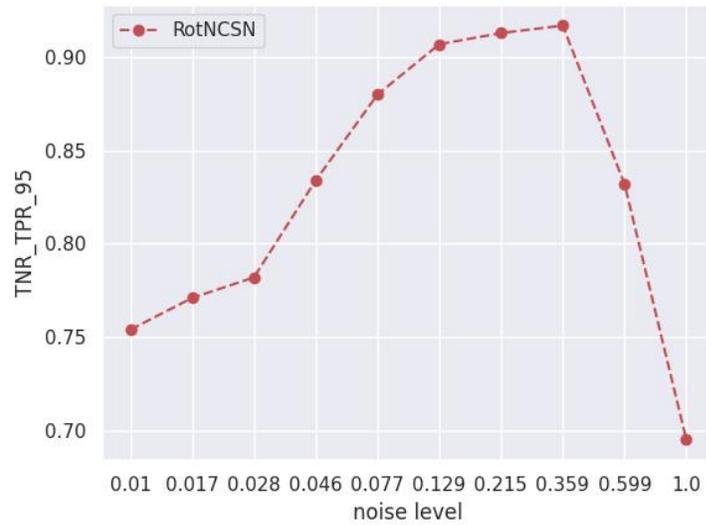


Figure 4: True negative rate (TNR) at 95% True positive rate (TPR) performance of RotNCSN w.r.t noise-level.