

Supplementary Material

1. On the influence of feature norm to L2 reconstruction error

In order to set the theoretical ground for supporting our normalized L2 distance introduced in main paper, we discuss and prove in this section the fact that feature with smaller norm leans to produce smaller reconstruction error.

Since the input of our autoencoder is one-dimensional AV feature, it is natural to employ FC network as the architecture of the encoder and decoder. Noticeably, the convolutional layer, shortcut connection and average pooling are all in essence linear mappings, and therefore the additional use of them will not concern the following discussion. For a FC autoencoder with L layers, we denote the weight matrix of l^{th} FC layer as $W^l \in \mathbb{R}^{n^l \times n^{l-1}}$, the offset bias as $\mathbf{b}^l \in \mathbb{R}^{n^l}$, the activation function as $\sigma(\cdot)$ and the input as $\mathbf{x} \in \mathbb{R}^H$. Then the pre-activation output of l^{th} layer could be recursively written as

$$f^l(\mathbf{x}) = W^l \sigma(f^{l-1}(\mathbf{x})) + \mathbf{b}^l. \quad (1)$$

Although activation functions in the context of deep learning are always nonlinear in the full space, they could be approximately considered as linear in a certain polytope. For instances, given an input in \mathbb{R}^1 , relu is linear if only the region of $[0, +\infty)$ or $(-\infty, 0)$ is considered. Sigmoid could be approximated as linear exclusively in the region of $(-\infty, -\alpha)$, $[-\alpha, \alpha]$ or $[\alpha, +\infty)$, for some $\alpha > 0$. For simplicity, in the following discussion we consider relu as the applied activation function $\sigma(\cdot)$. Similar to [2], our FC autoencoder can be expressed as a piecewise affine function

$$\begin{aligned} f^L(\mathbf{x}) &= W^L \sigma(W^{L-1} \sigma(\dots \sigma(W^1 \mathbf{x} + \mathbf{b}^1) \dots) + \mathbf{b}^{L-1}) + \mathbf{b}^L \\ &= W^L \Lambda^{L-1}(\mathbf{x}) (W^{L-1} \Lambda^{L-2}(\mathbf{x}) (\dots \Lambda^1(\mathbf{x}) (W^1 \mathbf{x} + \mathbf{b}^1) \dots) + \mathbf{b}^{L-1}) + \mathbf{b}^L \\ &= \Gamma \mathbf{x} + B, \end{aligned} \quad (2)$$

where $\Lambda^l(\mathbf{x}) \in \mathbb{R}^{n^l \times n^l}$, for $l = 1, \dots, L-1$ are diagonal matrices defined as

$$\Lambda^l(\mathbf{x}) = \begin{bmatrix} \mathbb{1}(f_1^l(\mathbf{x}) > 0) & & \\ & \dots & \\ & & \mathbb{1}(f_{n^l}^l(\mathbf{x}) > 0) \end{bmatrix}, \quad (3)$$

and $\Gamma \in \mathbb{R}^{H \times H}$ and $B \in \mathbb{R}^H$ are matrices defined as

$$\begin{aligned} \Gamma &= W^L \left(\prod_{i=1}^{L-1} \Lambda^{L-i}(\mathbf{x}) W^{L-i} \right), \\ B &= \sum_{i=1}^{L-1} \left(\prod_{k=1}^{L-i} W^{L+1-k} \Lambda^{L-k}(\mathbf{x}) \right) \mathbf{b}^i + \mathbf{b}^L. \end{aligned} \quad (4)$$

We can further have

$$\begin{aligned} \|\mathbf{x} - f^L(\mathbf{x})\| &= \|\mathbf{x} - \Gamma \mathbf{x} - B\| \\ &\leq \|\mathbf{x} - \Gamma \mathbf{x}\| + \|B\| \\ &\leq \|I - \Gamma\| \|\mathbf{x}\| + \|B\|. \end{aligned} \quad (5)$$

The number of total possible variants of Γ and B is $2^{\sum_{i=1}^{L-1} n^i}$, and the specific forms of Γ and B are determined by which polytope defined as intersection of $\sum_{i=1}^{L-1} n^i$ half spaces in \mathbb{R}^H the input \mathbf{x} is in. Therefore, given an input, an upper bound of its L2 reconstruction error is definite, and is approximately proportional to its norm. This result clearly supports our claim in the main paper that *feature with smaller norm tends to have smaller L1 or L2 reconstruction loss*. For instances, given a feature with norm close to zero, its reconstruction error is approximately $\|B\|$. While for a feature with norm $\rightarrow \infty$, its reconstruction error tends to be arbitrarily larger than $\|B\|$.

2. Architecture of encoder and decoders

There are one encoder and two decoders employed in our OOD detection module, namely E , D_1 and D_2 . The encoder has a FC layer with C filters (no additive bias term) and a following softmax function with temperature scaling, in which C is the number of ID classes.

Each decoder is a simple three-layer FC network with swish as activation function. Specifically, D_1 consists of sequential operations: FC layer with H filters and bias term \rightarrow swish nonlinearity \rightarrow FC layer with H filters and bias term \rightarrow swish nonlinearity \rightarrow FC layer with H filters and bias term, where H is the dimensionality of the AV feature

fed into the encoder. D_2 has the following operations: FC layer with 512 filters and bias term \rightarrow swish nonlinearity \rightarrow FC layer with 256 filters and bias term \rightarrow swish nonlinearity \rightarrow FC layer with C filters and bias term.

3. Model configurations in ablation study

Due to the limited pages in main paper, we did not detail the models applied in Ablation study (Table.2 in paper). Here we report them in order. For specific details and replication, please refer to the code submitted with this material.

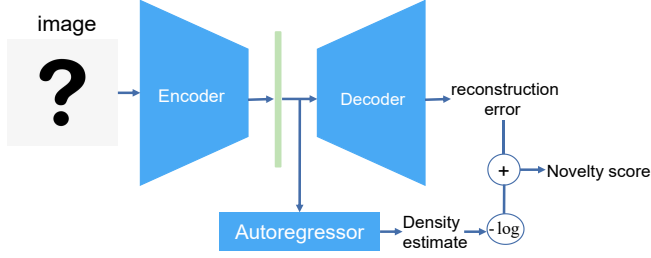


Figure 1. The overall framework of 1st LSA.

1st LSA: This is the original method of Latent Space Autoregression(LSA) [1]. Here we simply describe the framework of it, and a graphic representation is shown in Fig.1. With the training of a standard reconstruction autoencoder, this method trains an autoregressor on the latent codes output from the encoder to model the latent density while decreasing the differential entropy of the distribution of ID latent features. The input information of the autoencoder is the original image and the novelty score is computed as the summation of the L2 reconstruction error term and the negative logarithm of the density estimate produced from the autoregressor. Both terms are normalized using a validation set of ID data

$$novelty\ score = norm(L2(\mathbf{r}, \tilde{\mathbf{r}})) + norm(-\log(q(\mathbf{r}))),$$

in which \mathbf{r} denotes the input image, $\tilde{\mathbf{r}}$ the reconstruction of \mathbf{r} and $q(\mathbf{r})$ the density estimate of the latent code of \mathbf{r} . The model architecture and training settings are the same as those for CIFAR-10 dataset in the original paper.

2nd-image+feature: In this model, we change the input information of autoencoder from image to its AV feature extracted in the Wide-ResNet-28-10. Since the encoder and decoder in the original paper of LSA were designed to compress and rebuild image and therefore are parameter-redundant to be used in the case of low-dimensional AV features, we apply for this model a simpler self-designed autoencoder instead.

1. **Encoder:** FC layer with 100 filters and bias term \rightarrow relu nonlinearity \rightarrow FC layer with 100 filters and bias term \rightarrow sigmoid nonlinearity.

2. **Decoder:** FC layer with 640 filters and bias term \rightarrow swish nonlinearity \rightarrow FC layer with 640 filters and bias term \rightarrow swish nonlinearity \rightarrow FC layer with 640 filters and bias term.

3. **Autoregressor:** Masked FC layer with 32 channels \rightarrow leaky relu activation \rightarrow masked FC layer with 32 channels \rightarrow leaky relu activation \rightarrow masked FC layer with 32 channels \rightarrow leaky relu activation \rightarrow masked FC layer with 100 channels.

The applied novelty scoring function and other pipeline differ from 1st LSA only in that AV feature is considered rather than image.

3rd-L2+NL2: The framework of this model is identical to the above except that we use the proposed NL2 distance over the L2 distance as the $dist(\cdot)$ in the normality scoring function, which is calculated as

$$novelty\ score = norm(NL2(\mathbf{v}, \tilde{\mathbf{v}})) + norm(-\log(q(\mathbf{v}))),$$

where \mathbf{v} is the extracted AV feature.

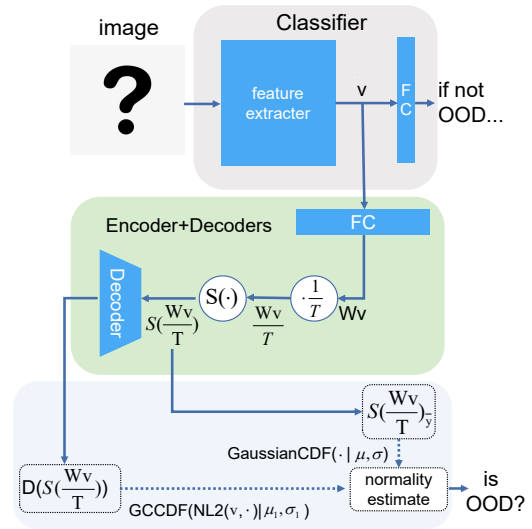


Figure 2. The overall framework of 4th -AutoReg+CE.

4th-AutoReg+CE: To restrict the span of latent representations, the models introduced above all apply a latent space autoregressor. In this model, we keep the decoder and change the encoder from that in 2nd-image+feature to the proposed one in order to use a simple cross entropy loss to compress the latent space. Thus, in test time the normality scoring function becomes

$$\begin{aligned}
\text{normality score} &= \Phi(S(\frac{W\mathbf{v}}{T})_{\bar{y}} | \mu, \sigma + 10\sigma) \\
&\cdot \Psi(\|\frac{\mathbf{v}}{\|\mathbf{v}\|} - \frac{D(S(\frac{W\mathbf{v}}{T}))}{\|\mathbf{v}\|} \| | \mu_1, \sigma_1 + 10\sigma_1).
\end{aligned}$$

The overall framework can be seen in Fig.2.

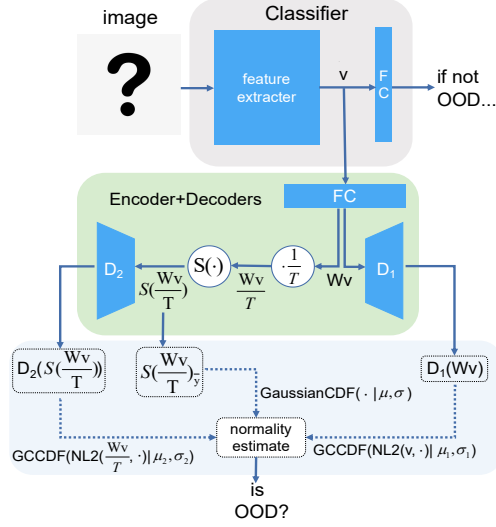


Figure 3. The overall framework of 5th-basic+layerwise.

5th-basic+layerwise: We apply the framework of layerwise reconstruction for this model. Thus it no longer recovers the input directly from the latent space and an additional decoder D_2 is involved. In essence, this model is identical to the one applied for benchmark experiments except that we do not use ϵ terms in the normality scoring function

$$\begin{aligned}
\text{normality score} &= \Phi(S(\frac{W\mathbf{v}}{T})_{\bar{y}} | \mu_0, \sigma_0) \\
&\cdot \Psi(\|\frac{\mathbf{v}}{\|\mathbf{v}\|} - \frac{D_1(W\mathbf{v})}{\|\mathbf{v}\|} \| | \mu_1, \sigma_1) \\
&\cdot \Psi(\|\frac{W\mathbf{v}}{\|W\mathbf{v}\|} - \frac{D_2(S(\frac{W\mathbf{v}}{T}))}{\|\frac{W\mathbf{v}}{T}\|} \| | \mu_2, \sigma_2).
\end{aligned}$$

The overall framework is presented in Fig.3.

6th+epsilon: This model is the one introduced and tested in benchmark experiment. Difference with respect to the model introduced above, we add the ϵ terms into the normality scoring function to prevent it from collapsing

$$\begin{aligned}
P(\mathbf{v} \in V) &= \Phi(S(\frac{W\mathbf{v}}{T})_{\bar{y}} | \mu_0, \sigma_0 + \epsilon_0) \\
&\cdot \Psi(\|\frac{\mathbf{v}}{\|\mathbf{v}\|} - \frac{D_1(W\mathbf{v})}{\|\mathbf{v}\|} \| | \mu_1, \sigma_1 + \epsilon_1) \\
&\cdot \Psi(\|\frac{W\mathbf{v}}{\|W\mathbf{v}\|} - \frac{D_2(S(\frac{W\mathbf{v}}{T}))}{\|\frac{W\mathbf{v}}{T}\|} \| | \mu_2, \sigma_2 + \epsilon_2).
\end{aligned}$$

Other configurations keep unchanged.

References

- [1] D. Abati, A. Porrello, S. Calderara, and R. Cucchiara. Latent space autoregression for novelty detection. 2018. [2](#)
- [2] Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 41–50, 2019. [1](#)