

Out of Distribution Detection Using Bayesian Variational Autoencoder

Kamran Chitsaz

Alireza Razaghi

1 Introduction

The recent success of deep learning methods plays an essential role in the evolution at an unprecedented rate of current technology. Despite their remarkable performance in a wide range of complex tasks including image recognition [1], object detection [2], and speech recognition [3], deep learning models can make confident, but meaningless, predictions on unobserved data [4]. To perform reliable decisions, Deep learning models need to be able to determine whether the input data is an anomaly or outlier and significantly different from the training data. Such data are called out-of-distribution (OoD) data. Several approaches have been proposed to detect OoD data such as spectral anomaly detection techniques which try to find lower dimensional embeddings of the input data where OoD and normal data are expected to be separated from each other. The Principal components analysis (PCA) [5] and Autoencoders [6] are among the most popular dimension reduction methods for detecting outliers.

Likelihood estimators are one of the most promising OoD detection approach which work for probabilistic generative models. Such models can evaluate the likelihood of input data, and if a generative model fits the training data distribution well enough, it should assign high likelihood to samples from the similar distribution and low likelihood to OoD samples. So to detect whether a data is OoD or not, by this approach first trains a density estimator parameterized by θ , and then classifies an input \mathbf{x}^* as OoD based on a threshold on the density of \mathbf{x}^* , i.e., if $p(\mathbf{x}^* | \theta) < \tau$ [7]. In this project we used variational autoencoder to calculate log likelihood and detect OoD.

2 Background

One of the most popular models for density estimation is the Variational Autoencoder [8] which successfully can approximate the distribution of training data. By this generative model we can learn a joint model $p(x, z)$ of the data x and some latent variables z . The idea is that each data X has a corresponding latent variable Z that data X is caused by this Z , so we can marginalize out Z .

$$p(x) = \int p(x|z)p(z)dz$$

The conditional distribution $p(x|z)$ is a Gaussian. So, we have a kind of mixture of infinitely many Gaussians, for each value of Z , there's one Gaussian and we mix them with weights. In this approach we assumed that prior is normal and likelihood is a Gaussian which it's parameters are depend on latent variable z . One way to define parameters is using neural

networks. So the parameters of our model depends on leading variable z which is generated by a neural network with parameter θ . So we can build the observation model as follow:

$$\begin{aligned} p(x) &= \int p(x|z)p(z)dz \\ p(z) &= \mathcal{N}(0, I) \\ p(x | z) &= \mathcal{N}(x | \mu_\theta(z), \text{diag}(\sigma_\theta^2(z))) \end{aligned}$$

Since computing $p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x} | \mathbf{z})p(\mathbf{z})d\mathbf{z}$ is intractable in high dimension, variational inference is used to derive a lower bound on the log likelihood of \mathbf{x} . So we can approximate posterior $p_\theta(\mathbf{z} | \mathbf{x})$ using amortized variational inference via another neural network with parameter ϕ . Similarly, we can build the inference model as

$$q(z | x) = \mathcal{N}(z | \mu_\phi(x), \text{diag}(\sigma_\phi^2(x)))$$

Using Jensen's inequality we have

$$\begin{aligned} \log p(x) &= \log \int \frac{q(z | x)}{q(z | x)} p(x, z) dz \geq \int q(z | x) \log \frac{p(x, z)}{q(z | x)} dz \\ &= E_{q(z|x)} \log \frac{p(x, z)}{q(z | x)} = E_{q(z|x)} \log \frac{p(x | z)p(z)}{q(z | x)} \end{aligned}$$

So the objective we're trying to maximize is

$$\begin{aligned} \underset{w, \phi}{\text{maximize}} \quad & \sum_i E_{q_i} \log \frac{p(x_i | z_i, w)p(z_i)}{q_i(z_i)} \\ \text{subject to} \quad & q_i(t_i) = \mathcal{N}(m(x_i, \phi), \text{diag}(s^2(x_i, \phi))) \end{aligned} \tag{1}$$

We can rewrite equation 1 as below:

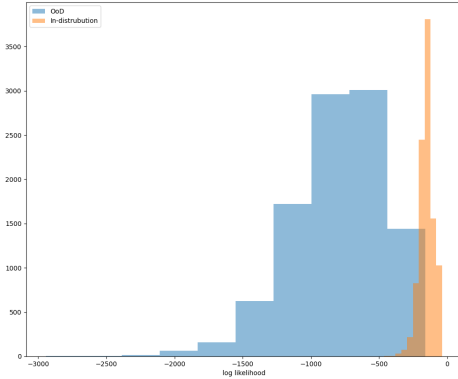
$$\begin{aligned} & \sum_i E_{q_i} \log p(x_i | t_i, w) + E_{q_i} \log \frac{p(t_i)}{q_i(t_i)} \\ &= \sum_i E_{q_i} \log p(x_i | t_i, w) - \mathcal{KL}(q_i(t_i) \| p(t_i)) \end{aligned} \tag{2}$$

When we maximize this minus KL we are actually trying to minimize KL so we are trying to push the variational distribution q_i as close to the prior as possible which is standard normal. Also if for simplicity we set all the output variances to be 1, the first term, log likelihood of x_i given t_i , can be interpreted as $-\|x_i - \mu(t_i)\|^2 + \text{const}$ which is actually a reconstruction loss. It tries to push x_i as close to the reconstruction as possible.

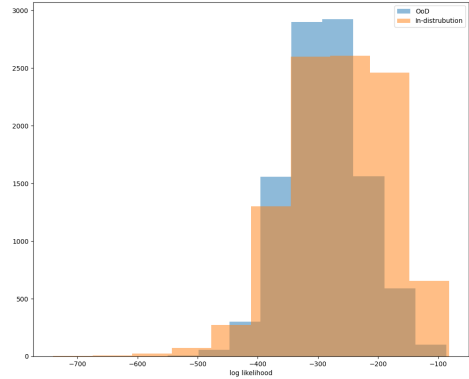
Finally, one can use importance sampling w.r.t. the variational posterior $q(z | x, \phi)$ to get an estimator $\hat{p}(\mathbf{x} | \theta, \phi)$ of the probability $p(\mathbf{x} | \theta)$ of an input \mathbf{x} under the generative model, i.e.,

$$p(\mathbf{x} | \theta) \simeq \hat{p}(\mathbf{x} | \theta, \phi) = \frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x} | \mathbf{z}_k, \theta) p(\mathbf{z}_k)}{q(\mathbf{z}_k | \mathbf{x}, \phi)}$$

where $\mathbf{z}_k \sim q(\mathbf{z} | \mathbf{x}, \phi)$ and where the estimator $\hat{p}(\mathbf{x} | \theta, \phi)$ is conditioned on both θ and ϕ to make explicit the dependence on the parameters ϕ of the proposal distribution $q(z | x, \phi)$.



(a)



(b)

Figure 1: Histogram of the log likelihood of test samples for (a) VAE trained on Mnist and Fashion-Mnist as OoD, and (b) VAE trained on Fashion-Mnist and Mnist as OoD

3 Results

We have trained our VAEs with samples only from training data from in distribution dataset and evaluate our method by test sample from different dataset to measure OoD performances. In first experiment we have set Mnist dataset as in distribution and Fashion-Mnist as OoD which you the results are illustrated in Fig. 1a. In this case we observed that our method assign higher probability to in distribution data and can detect OoD with accuracy 89.90%. But as another experiment we have set Fashion-Mnist dataset as in distribution and Mnist as OoD which you the results are illustrated in Fig. 1b. As you can see our method failed. Also, reference [9] shows that almost all major types of probabilistic generative models, including VAE, can assign spuriously high likelihood to OoD samples which make them unreliable metrics for OoD detection.

3.1 Metrics

To classify whether a data is OoD or not we compute maximum likelihood and then compare with a threshold. The choice of a threshold depends on the particular application. We have tune the threshold, as same as other model hyper parameters, by validation set.

References

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [2] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- [3] Li Deng, Jinyu Li, Jui-Ting Huang, Kaisheng Yao, Dong Yu, Frank Seide, Michael Seltzer, Geoff Zweig, Xiaodong He, Jason Williams, et al. Recent advances in deep learning for speech research at microsoft. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 8604–8608. IEEE, 2013.
- [4] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 427–436, 2015.
- [5] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3):1–58, 2009.
- [6] Mayu Sakurada and Takehisa Yairi. Anomaly detection using autoencoders with non-linear dimensionality reduction. In *Proceedings of the MLSDA 2014 2nd Workshop on Machine Learning for Sensory Data Analysis*, pages 4–11, 2014.
- [7] Christopher M Bishop. Novelty detection and neural network validation. *IEEE Proceedings-Vision, Image and Signal processing*, 141(4):217–222, 1994.
- [8] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [9] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do deep generative models know what they don’t know? *arXiv preprint arXiv:1810.09136*, 2018.