
Interpreting Rate-Distortion of Variational Autoencoder and Using Model Uncertainty for Anomaly Detection

Seonho Park George Adosoglou Panos M. Pardalos

Department of Industrial and Systems Engineering

University of Florida

Gainesville, Florida, USA

{seonhopark, g.adosoglou, pardalos}@ufl.edu

Abstract

Building a scalable machine learning system for unsupervised anomaly detection via representation learning is highly desirable. One of the prevalent methods is using a reconstruction error from variational autoencoder (VAE) via maximizing the evidence lower bound. We revisit VAE from the perspective of information theory to provide some theoretical foundations on using the reconstruction error, and finally arrive at a simpler and more effective model for anomaly detection. In addition, to enhance the effectiveness of detecting anomalies, we incorporate a practical model uncertainty measure into the metric. We show empirically the competitive performance of our approach on benchmark datasets.

1 Introduction

Autoencoders have been widely used in many machine learning applications not only to reduce the noise from the input to learn representations but also to reconstruct the output with the salient information of the input. These autoencoders learn common information of the inputs by mapping to the latent representations in an unsupervised manner. When it comes to anomaly detection, using the reconstruction error of various autoencoders to discern anomalies has been widely and successfully employed [1, 2, 3, 4, 5], even though using reconstruction error lacks its theoretical foundations. One of the autoencoders whose theoretical basis comes from variational inference is variational autoencoders (VAEs) [6]. VAEs try to minimize the difference between the true posterior and the variational posterior via maximizing the evidence lower bound (ELBO) with respect to the neural networks based encoder and decoder. After training, we expect that the ELBO approximates the marginal likelihood of the data.

In this work, we revisit VAEs from the perspective of rate-distortion theory [7, 8] to elucidate the roles of the two terms: the rate and distortion. Also, we argue that various autoencoders including β -VAE [9] can be explained with the trade-off between the rate and distortion in this perspective. Then, for the purpose of anomaly detection, we show that using only the encoder is more efficient to approximate the marginal likelihood, and finally, we arrive at a much simpler and more efficient model to discern anomalies.

Moreover, in order to enhance the performance of detecting anomalies, we incorporate the model uncertainty into our anomaly detection score. Since anomalies are unseen when the model is trained in an unsupervised setting, model uncertainty can capture the anomalies for which the model's

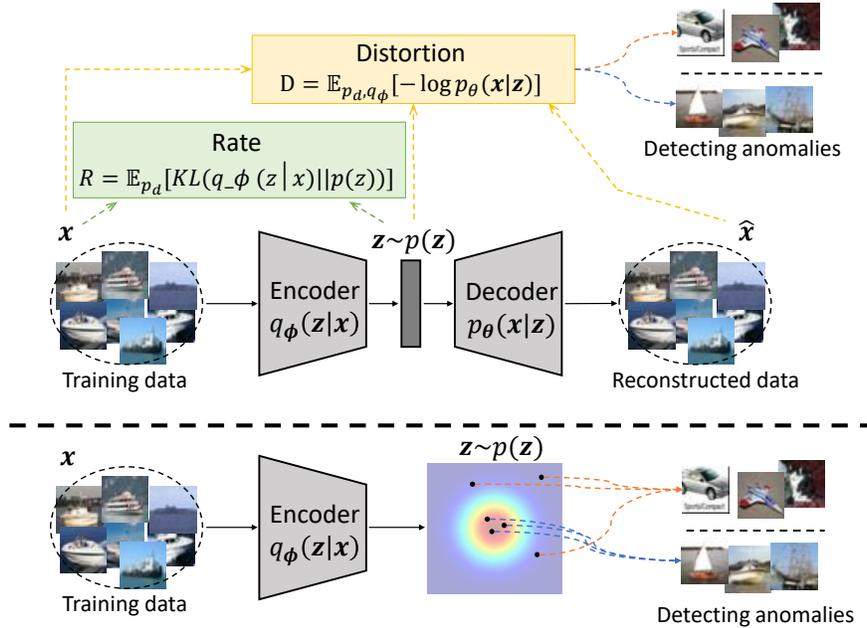


Figure 1: Overview of the proposed approach for anomaly detection. *Top*: it shows the VAE based anomaly detection approach from the perspective of rate-distortion theory. The reconstruction error (distortion measure) is used as a metric of anomaly detection. *Bottom*: we only use the encoder (without the decoder) of VAEs to identify anomalies.

confidence is low. Finally, we justify our approaches with theoretical analyses as well as experiments on benchmark datasets.¹

To summarize, we make the following contributions for anomaly detection.

- With the theoretical foundation on variational inference and rate-distortion theory, we elucidate that the VAE-based anomaly detection system aim to approximate the marginal probability of the data.
- We propose that using the encoder only is more efficient and simpler than VAE’s reconstruction error to measure anomaly score.
- We incorporate model uncertainty into the metric to enhance anomaly detection performance.
- We provide theoretical and empirical basis on our approach for anomaly detection.

2 Problem Definition

What we aim to do in this work is to derive an anomaly score $s(\mathbf{x})$ to indicate whether a given query datapoint \mathbf{x} is anomalous or not. More formally, with a scalar threshold γ , anomaly score $s(\mathbf{x})$ should distinguish anomalous instances as,

$$s(\mathbf{x}) \geq \gamma \Rightarrow \text{anomalous}$$

$$s(\mathbf{x}) < \gamma \Rightarrow \text{normal}$$

We cannot anticipate which anomalies come to the system in the future so it is reasonable to assume that an anomaly detection model is learned in an unsupervised way, that is, when training, we only have access to normal data and when testing we can access *contaminated* data consisting of both normal and anomalous instances. This setting is also referred to as one class classification [10, 11].

¹One can reach out to the public implementation for whole experiments via https://github.com/seonho-park/PNG_anomaly_detection

3 Information Theoretical Interpretation of VAE

In this section, we revisit VAEs [6] in the context of information theory to clarify the terms of the loss function of VAEs.

Variational Autoencoder Let us assume that we have a dataset $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ consisting of normal datapoints $\mathbf{x} \in \mathcal{X}$, i.i.d. sampled. The datapoints in \mathbf{X} are realized by a random process, $p^*(\mathbf{x}|\mathbf{z})p^*(\mathbf{z})$, where $p^*(\mathbf{z})$ and $p^*(\mathbf{x}|\mathbf{z})$ are a true prior over latent variables and a true likelihood, respectively. Also, we assume that the latent random variable $\mathbf{z} \in \mathcal{Z} \subseteq \mathbb{R}^J$ follows a true prior $p^*(\mathbf{z})$.

Given an input \mathbf{x} , a variational posterior (which is also referred to as an encoder) is derived to approximate a true posterior via KL divergence and the corresponding marginal log-likelihood can be expressed with variational inference (VI) as

$$\log p(\mathbf{x}) = KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})) + \mathcal{L}_{VI}(\phi, \theta; \mathbf{x}) \quad (1)$$

where $q_\phi(\mathbf{z}|\mathbf{x})$ is a neural network model parameterized by parameters ϕ . The first RHS term of the above equation (Eq. 1) is the KL divergence between the variational and true posterior. The KL divergence is always nonnegative and it is zero if and only if the variational posterior is exactly equivalent to the true posterior which is intractable to compute directly. Because the KL divergence is nonnegative, we could say the second RHS term is the lower bound of the marginal log-likelihood, $\log p(\mathbf{x})$ which is fixed. This second RHS term can be elaborated as

$$\mathcal{L}_{VI}(\phi, \theta; \mathbf{x}) = -KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] \quad (2)$$

where $p_\theta(\mathbf{x}|\mathbf{z})$ is a variational approximation (decoder) to a true likelihood, parameterized by parameters θ and $p(z)$ is an approximation to $p^*(z)$. The first RHS term acts as a regularizer of q_ϕ and the second RHS term corresponds to the negative reconstruction error. By taking an expectation w.r.t. the empirical data distribution $p_d(\mathbf{x})$, VAEs seek to maximize the evidence lower bound (ELBO) to minimize the KL divergence between the variational posterior and true posterior as

$$\begin{aligned} \max_{\phi, \theta} \mathbb{E}_{p_d(\mathbf{x})} [-KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))] \\ + \mathbb{E}_{p_d(\mathbf{x})} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]] \end{aligned} \quad (3)$$

The ELBO consists of two terms. The first term in Eq.3 can be interpreted as compression loss of the input information. If the first term is high (as $q_\phi(\mathbf{z}|\mathbf{x})$ approaches $p(\mathbf{z})$), it means that the latent code compresses the input so well that the salient information of the input disappears.² The second term is the expected negative reconstruction error, which represents the (negative) difference between the input and the output from the decoder. Thus, the ELBO can be interpreted as the trade-off between the compression loss (how much information can be lost in the latent space) and the (negative) reconstruction error (how much information can be retrieved from the decoder).

VAE as Lossy Compression From the perspective of the rate-distortion theory [7], we revisit VAE to elucidate the roles of the terms of the ELBO. We derive two terms the *rate* and *distortion*, which correspond to negative compression loss and the reconstruction error, respectively.

Based on the previous work [8], we can rewrite the VAE problem as:

$$\min_{\phi, \theta} R(\mathbf{x}, \mathbf{z}|\phi) + D(\mathbf{x}, \mathbf{z}|\phi, \theta) \quad (4)$$

where

$$R(\mathbf{x}, \mathbf{z}|\phi) = \mathbb{E}_{p_d(\mathbf{x})} [KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))] \quad (5)$$

$$D(\mathbf{x}, \mathbf{z}|\phi, \theta) = \mathbb{E}_{p_d(\mathbf{x})} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [-\log p_\theta(\mathbf{x}|\mathbf{z})]] \quad (6)$$

The rate, R , is the expected value of the rate measure, $KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$. The rate represents the expectation of the KL divergence between the encoder and prior.

²Sometimes, compression is also referred to as disentanglement because what we aim to get as a latent representation is usually a disentangled representation and manipulate some elements of the latent vector to tweak the reconstruction readily. Please see [12, 13, 9] for more details.

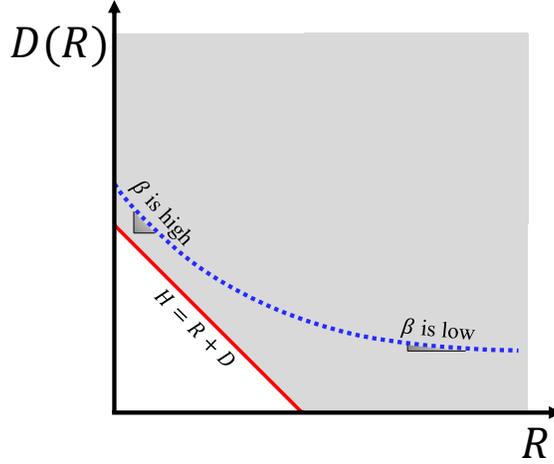


Figure 2: Schematic view of distortion-rate function. A red line corresponds to the theoretical lower bound of the rate and distortion. By varying β of β -VAE, we could achieve the points on a blue dashed curve, the sub-optimal distortion-rate function, which is best achievable with VAEs.

D is the distortion, the expected value of distortion measure, $d(\mathbf{x}, \mathbf{z}) = \mathbb{E}_{q_\phi}[-\log p_\theta]$ representing the reconstruction error. Note that the rate only depends on the parameters ϕ of the encoder, while the distortion depends on both ϕ and θ .

Also, we would like to introduce the data entropy, H , as

$$H(\mathbf{x}) = \mathbb{E}_{p_d(\mathbf{x})}[-\log p(\mathbf{x})] \quad (7)$$

Given R , D , and H , the expectation w.r.t. $p_d(\mathbf{x})$ of Eq.1 can be rewritten as,

$$H = \mathbb{E}_{p_d(\mathbf{x})}[-KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x}))] + R + D \quad (8)$$

From the nonnegative property of the KL divergence, we can say $H \leq R + D$ where the equality holds if and only if the variational posterior equals to the true posterior, i.e., $q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = p(\mathbf{z}|\mathbf{x}^{(i)})$, $\forall \mathbf{x}^{(i)} \in \mathbf{X}$. Then, $H = R + D$. This represents a theoretical lower bound of $R + D$ and is depicted as the red solid line in Fig. 2. In VAE, we may not achieve this ideal case, $H = R + D$, because of the limited finite families of parameters, the approximated prior and noises in the given dataset. Instead, we seek to find the (information) distortion-rate function (curve) by solving the following optimization problem:

$$\begin{aligned} & \min_{\phi, \theta} D \\ & \text{subject to } R \leq \bar{R} \end{aligned} \quad (9)$$

where \bar{R} denotes an upper limit of the rate. In order to optimize both without taking \bar{R} , we can take the Lagrangian of Eq. 9 with a Lagrange multiplier $\beta > 0$ as³,

$$\min_{\phi, \theta} D + \beta R \quad (10)$$

which resembles the β -VAE objective [9].

Varying β in β -VAE, we can get the distortion-rate function depicted as a blue dashed curve in Fig. 2. Even though the curve is not explicitly formed except for some known simple examples, it is known that the distortion-rate function is convex and monotonically non-increasing. With these properties, the Lagrangian multiplier can be interpreted as a negative slope of the distortion-rate function. This β plays a role to balance the rate and the distortion. When β is high, we can get the point where the rate is low and the distortion is high. Whereas, when β is low, we can get the point where the rate is high and the distortion is low. Because the joint distribution $p(\mathbf{x}, \hat{\mathbf{x}})$ of lossy compression is composed of the encoder and decoder, it is noted that distortion cannot be zero even when the rate is high enough [14].

³One could think this formulation is to minimize two objectives R and D . Then, we can achieve the Pareto frontier of the objectives, which is corresponding to the distortion-rate function in rate-distortion theory as well.

4 Method

Zero Rate Setting to Approximate Data Distribution In this section, we first introduce a method to approximate the marginal log-likelihood of the data and derive the anomaly score considering model uncertainty as well as the approximate data distribution to detect anomalies.

The main purpose of using VAEs is to encode the salient information of the data to the latent space and reconstruct the output $\hat{\mathbf{x}}$. When it comes to anomaly detection, even though using the reconstruction error of autoencoders shows great performances empirically, it is not so straightforward and lacks some theoretical foundations. From Eq. 9, if we set R to be zero then the learned rate has to be zero. $R = 0$ means from the definition of the rate that $q_\phi(\mathbf{z}|\mathbf{x}) = p(\mathbf{z})$ for all datapoints in training dataset. The latent variable \mathbf{z} does not store any particular information of the individual datapoint and the decoder seeks to give outputs via the stochastic decoder and the reparameterization trick [6] which resemble the empirical data distribution, $p_d(\mathbf{x})$. This means that the decoded outputs of different inputs are widely distributed as the training datapoints are. Precisely, the zero rate corresponding distortion can be rewritten with the definition of the distortion (Eq. 6) as

$$D = - \int d\mathbf{x} p_d(\mathbf{x}) \int d\mathbf{z} p(\mathbf{z}) \log p_\theta(\mathbf{x}|\mathbf{z}) \quad (11)$$

Thus, the gap between H and D when $R = 0$ can be decreased when we use sufficiently large $p_d(\mathbf{x})$ and a powerful decoder p_θ . Equivalently, the difference of intercepts of the blue dashed curve and the red line in Fig. 2 can be interpreted by means of insufficient data representations of dataset, poor decoder performance, and implicitly wrong selection of prior and its dimension.

From this observation, we could conclude that, for the purpose of anomaly detection, we can achieve the distortion-rate function by setting $\beta = \infty$ as well. Then, we only need the encoder of the autoencoder to estimate the marginal log-likelihood of the data. As a consequence, we solve the following problem to estimate the marginal log-likelihood of the data:

$$\min_{\phi} \mathcal{L}(\phi) = \mathbb{E}_{p_d(\mathbf{x})} [KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))] \quad (12)$$

With this formulation, since the decoder p_θ vanishes, we could have a more concise model (only with the encoder) rather than using the reconstruction error which requires both the encoder and the decoder. The model can be learned with first or second order stochastic optimization methods such as stochastic gradient descent, ADAM [15], AdaGrad [16], Newton or its variant methods [17, 18] via direct backpropagation.

Prior Generating Networks We can break down the rate as reported in [19] as

$$\mathbb{E}_{p_d} [KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))] = I(\mathbf{x}; \mathbf{z}) + KL(q_\phi(\mathbf{z})||p(\mathbf{z})) \quad (13)$$

where I denotes the mutual information between \mathbf{x} and \mathbf{z} under the joint distribution $q_\phi(\mathbf{x}, \mathbf{z}) = q_\phi(\mathbf{z}|\mathbf{x})p_d(\mathbf{x})$. Also, $q_\phi(\mathbf{z})$ is known as the aggregated variational posterior [20, 13] and can be attained as $q_\phi(\mathbf{z}) = \mathbb{E}_{p_d} [q_\phi(\mathbf{z}|\mathbf{x})]$. Also, when we set $p(\mathbf{z}) = q_\phi(\mathbf{z})$, the rate is identical to the mutual information $I(\mathbf{x}; \mathbf{z})$. This setting of the prior is referred to as the VampPrior [20]. The mutual information $I(\mathbf{x}; \mathbf{z})$ is upper bounded by the rate R as above. Therefore, the zero rate setting means that the latent codes do not have any information about normal datapoints in the training dataset. Intuitively, the rate measure $KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$ of anomalous instance is higher than that of normal instances.

In our experiments, the prior $p(\mathbf{z})$ is defined by an isotropic multivariate Gaussian, $\mathcal{N}(\mathbf{0}, \mathbf{1})$ as in [6]. Also let the encoder be the neural network based model of which the outputs are the mean and standard deviation of the isotropic multivariate Gaussian, i.e., $q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = (\boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{(i)})$. Thus $\mathcal{L}(\phi)$ in Eq.12 can be simplified as

$$\begin{aligned} & \mathbb{E}_{p_d(\mathbf{x})} [KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))] \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \left(\log \frac{1}{\sigma_j^{(i)}} + \frac{(\sigma_j^{(i)})^2 + \mu_j^{(i)^2}}{2} - \frac{1}{2} \right) \end{aligned} \quad (14)$$

In our experiments, we have used this setting and we call this *prior generating networks (PGN)* in what follows. This is named after the fact that the neural network based encoder merely aims to approximate the prior $p(\mathbf{z})$.

Anomaly Score with Model Uncertainty Disregarding the expectation with respect to p_d , from Eq. 8 we can estimate the log probability of a query input \mathbf{x}^* . Note that as the encoder converges to the prior we cannot reconstruct the input, the distortion measure should be constant. Hopefully, if we assume that the KL divergence between the variational posterior and true posterior is sufficiently low, so negligible, then the log probability of \mathbf{x}^* is proportional to the negative KL divergence between the variational posterior and prior, i.e., $\log p(\mathbf{x}^*) \propto -KL(q_\phi(\mathbf{z}|\mathbf{x}^*)||p(\mathbf{z}))$.

Our model $q_\phi(\mathbf{z}|\mathbf{x})$ can be deterministic and the data distribution, $p_d(\mathbf{x})$, can only impose some stochasticity into the model. The model can gain more stochasticity by employing a random noise ξ into the model as $q_\phi(\mathbf{z}|\mathbf{x}, \xi)$. One practical way to do this is to use MC dropout [21]. Inserting dropout layers [22] in the model, MC dropout estimates the first and second moments by Monte Carlo samplings with T stochastic forward passes.

MC dropout is one of the most prevalent methods used for capturing model uncertainty. Model uncertainty, also referred to as *epistemic* uncertainty, comes from the lack of knowledge of the data. It includes uncertainties generated by the situation where the model does not have enough knowledge and/or experience on the data⁴. From the problem definition, we also assume that the model is trained only on the normal datapoints so that the model uncertainty can capture anomalies by generating higher uncertainties on them.

Let us define $\xi^{(t)}$ as the t -th realization of the random noise ξ . Also, its elements are i.i.d and sampled from the Bernoulli distribution with the dropout probability p . With the T stochastic forward passes, $KL(q_{\phi^*}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) = KL(\mathbb{E}_\xi [q_{\phi^*}(\mathbf{z}|\mathbf{x}, \xi)] ||p(\mathbf{z}))$ and $\mathbb{E}_\xi [q_{\phi^*}(\mathbf{z}|\mathbf{x}, \xi)] \simeq \frac{1}{T} \sum_{t=1}^T q_{\phi^*}(\mathbf{z}|\mathbf{x}, \xi^{(t)})$. For simplicity, we only impose ξ into the mean output μ of q_ϕ . As a result, we propose and use the following anomaly score metric. Given a query point \mathbf{x}^* and learned parameters ϕ^* ,

$$s(\mathbf{x}^*) = \frac{1}{T} \sum_{t=1}^T KL \left(q_{\phi^*}(\mathbf{z}|\mathbf{x}^*, \xi^{(t)}) ||p(\mathbf{z}) \right) = \quad (15)$$

$$KL \left(\frac{1}{T} \sum_{t=1}^T q_{\phi^*}(\mathbf{z}|\mathbf{x}, \xi^{(t)}) ||p(\mathbf{z}) \right) + \text{Variation}[\mu(\mathbf{x}^*)] \quad (16)$$

where $\mu(\mathbf{x}^*)$ is an abbreviation of the mean output of $q_\phi(\mathbf{z}|\mathbf{x}^*)$ and $\text{Variation}[\cdot]$ denotes the model uncertainty measure. It is noted that Eq. 15 is equivalent to the summation of MC dropout based estimation and the measured model uncertainty. So Eq. 15 as an anomaly score is efficient to capture both the mean KL value with additionally imposed stochasticity and model uncertainty via MC dropout.

Theoretical Analysis of Model Uncertainty We elucidate the reason that Eq. 15 is equal to Eq. 16. To this end, we first revisit the following lemma regarding the gap of Jensen’s inequality.

Lemma 1 (The Gap of Jensen’s inequality). *Let x be a one dimensional random variable and $p(x \in (a, b)) = 1$ where $-\infty \leq a < b \leq \infty$. Let $\varphi(x)$ be a twice differentiable function on (a, b) . Then,*

$$\begin{aligned} \frac{\inf_x \varphi''(x) \text{Var}(x)}{2} &\leq \mathbb{E}[\varphi(x)] - \varphi(\mathbb{E}[x]) \\ &\leq \frac{\sup_x \varphi''(x) \text{Var}(x)}{2} \end{aligned} \quad (17)$$

Proof. Please refer to the proof of the theorem 1 in [24] and [25]. □

This lemma implies that when the function φ is strictly convex the Jensen’s inequality gap represents the variance of the random variable. With this, we can derive the following theorem to justify Eq. 15, the proposed anomaly score.

⁴The protocol of OOD detection is similar to that of anomaly detection, where we have access to normal data when training and distinguish the anomalies (or data came from other datasets). However, OOD detection is usually conducted as a byproduct to assess the “confidence” of the system for classifications or regressions while the system for anomaly detection is merely for it. Please see [23] for more details on OOD detection.

Theorem 2 (PGN anomaly score measure). *From a finite parameters set ϕ , let us assume that we have learned parameters ϕ^* . Let $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ be outputs of the encoder q_ϕ and $p(\mathbf{z})$ be an isotropic multivariate Gaussian prior. Also, assume that $\boldsymbol{\mu}$ involves a random noise $\boldsymbol{\xi}$. Then, given an arbitrary input data \mathbf{x}^* , the following equality holds*

$$\begin{aligned} & \mathbb{E}_{\boldsymbol{\xi}} [KL(q_{\phi^*}(\mathbf{z}|\mathbf{x}^*, \boldsymbol{\xi})||p(\mathbf{z})))] = \\ & KL(\mathbb{E}_{\boldsymbol{\xi}} [q_{\phi^*}(\mathbf{z}|\mathbf{x}^*, \boldsymbol{\xi})] ||p(\mathbf{z})) + \alpha \sum_{j=1}^J \text{Var}_{\boldsymbol{\xi}}(\mu_j(\mathbf{x}^*|\boldsymbol{\xi})) \end{aligned} \quad (18)$$

Proof. Let us define $\varphi(\mu_j) := KL(\mathcal{N}(\mu_j, \sigma_j)||p(z_j))$ where $p(z_j) = \mathcal{N}(0, 1)$. From the fact that q_ϕ is finite valued and twice differentiable function based on the neural networks, KL divergence with a fixed $p(z_j)$ is also twice differentiable and strongly convex with respect to μ_j , i.e., $\varphi''(\mu_j) > 0, \forall \mu_j$. Let us denote that $\inf \varphi''(\mu_j) = m_j$ and $\sup \varphi''(\mu_j) = M_j$ where $0 < m_j < M_j < \infty \forall j \in \{1, J\}$. Because of Lemma 1, the following equality holds:

$$\frac{m_j}{2} \text{Var}(\mu_j) \leq \mathbb{E}[\varphi(\mu_j)] - \varphi(\mathbb{E}[\mu_j]) \leq \frac{M_j}{2} \text{Var}(\mu_j) \quad (19)$$

Summing Eq. 19 upto J and taking α such that $\min_j \{m_j\} < \frac{2\alpha}{J} < \max_j \{M_j\}$ finalize the proof. \square

From Eq. 16, model uncertainty metric, $\text{Variation}[\cdot]$, is proportional to $\sum_{j=1}^J \text{Var}_{\boldsymbol{\xi}}(\mu_j(\mathbf{x}^*|\boldsymbol{\xi}))$ of Eq. 18. The expected value can be approximated by T stochastic forward passes. Because it only needs T inferences and does not involve any further computations, incorporating model uncertainty into our anomaly score metric (Eq. 15) is so practical.

Relationship to Deep SVDD We would like to highlight that our PGN learning (Please see Eq. 12) generalizes the popular anomaly detection method, Deep SVDD [11]. Deep SVDD is to minimize the hypersphere in the latent space. If an instance lies out of the hypersphere, it is deemed anomalous. Given an input dataset \mathbf{X} and a prescribed center $\mathbf{c} \in \mathcal{Z}$, *One-Class Deep SVDD* [11] trains the neural network based model $f_{\mathbf{w}}$, parameterized by parameters \mathbf{w} , as

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N \left\| f_{\mathbf{w}}(\mathbf{z}|\mathbf{x}^{(i)}) - \mathbf{c} \right\|^2 + \lambda \Omega(\mathbf{w}) \quad (20)$$

where the second term represents a weight decay regularizer with a hyperparameter $\lambda > 0$. After training, the anomaly score can be calculated as

$$s(\mathbf{x}) = \|f_{\mathbf{w}^*}(\mathbf{z}|\mathbf{x}) - \mathbf{c}\|^2 \quad (21)$$

where \mathbf{w}^* are the learned parameters.

From Eq.12, suppose that we set the prior $p(\mathbf{z})$ to an isotropic multivariate Gaussian, $\mathcal{N}(\mathbf{c}, \epsilon \mathbf{I})$, with $\epsilon \ll 1$ and that q_ϕ only gives the estimated mean of the Gaussian with a fixed variance, i.e., $q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}^{(i)}, \epsilon \mathbf{I})$. Then, Eq.12 can be rewritten as

$$\begin{aligned} & \min_{\phi} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \left(\log \frac{\epsilon}{\epsilon} + \frac{\epsilon + (\mu_j^{(i)} - c_j)^2}{2\epsilon} - \frac{1}{2} \right) \\ & = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J \left(\frac{(\mu_j^{(i)} - c_j)^2}{2\epsilon} \right) \end{aligned} \quad (22)$$

Therefore, disregarding the weight decay term in Eq.20, we claim that PGN is, in some sense, a general formulation of the Deep SVDD and it provides different perspective to Deep SVDD, which does not rely on the previous kernel based methods such as OC-SVM [10] or SVDD [26].

5 Related Works

Deep Anomaly Detection based on Neural Networks Outlier detection using replicator neural networks [1] is, to the best of our knowledge, the first anomaly detection that uses neural networks where the reconstruction error is used as an anomaly score named ‘*outlyingness score*’. They introduced the replicator neural networks, feed-forwarding multi-layer perceptron neural networks with three hidden layers that forms the compressed latent representations and tries to reconstruct the inputs. Many recent approaches that are based on the reconstruction error of autoencoders [2, 3, 27] are also based on the same philosophical reasons. To enhance the performance of detecting anomalies in a huge amount of complex and high-dimensional data, the different types of deep autoencoders such as VAEs [6], adversarial autoencoders (AAEs) [28], denoising autoencoder [29], and deep convolutional autoencoders [30] have been equipped.

In a similar vein, generative adversarial networks (GANs) [31] have been also used as architectures of anomaly detection while using the reconstruction error as an anomaly score. The examples of this approach include AnoGAN [32] and OCGAN [33]. Because GANs focus on the powerful data generation, these anomaly detection approaches suggest some methodologies to increase the reconstruction error for anomalies. There are also deep neural networks based anomaly detection methods by *One-class classification* such as Deep SVDD [11] and OC-NN [34], which are inspired by kernel based methods, SVDD [26] and OC-SVM [10], respectively.

Generative Probabilistic Novelty Detection (GPND) [35] seeks to approximate the probability density of the data to distinguish anomalies. To this end, they train the AAE-like architecture to learn the manifold structure of data distribution.

Understanding VAE with Information Theory *Information bottleneck* [36] can be utilized to understand representation learning such as variational autoencoders (VAEs). VAEs is regarded as lossy compression in the context of rate-distortion theory. This way of understanding gives a different theoretical point of view to understand VAEs. [8, 37, 38]. In [39], the author argues that considering rate-distortion theory is the key to understanding representation learning and studies the possibility to optimize the marginal prior which is usually treated as fixed. Alemi et al. [8] try to understand ELBO with the rate-distortion theory based framework which is similar to our work. They show that decoupling ELBO to the rate and distortion helps understand the behavior of VAEs.

6 Experiments

6.1 Distortion-Rate Functions of VAE

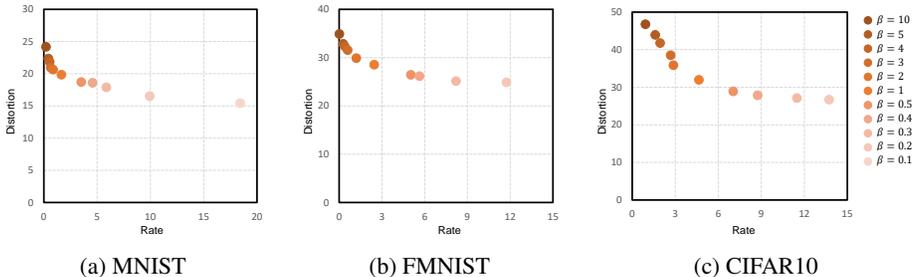


Figure 3: Distortion-rate functions on MNIST, FMNIST and CIFAR10 datasets. These show that VAE architectures which we use for anomaly detection experiments have sub-optimal distortion-rate functions empirically.

Experiment Settings To verify our assumption that VAEs have the sub-optimal distortion-rate functions as depicted in Fig. 2, we investigate the rate and distortion by varying β of β -VAE (Eq. 9) on MNIST [40], Fashion MNIST (FMNIST) [41] and CIFAR10 [42] datasets. We have used the official pre-split training dataset to train the model without any label information for all datasets. This work was similarly conducted by previous works [8, 37]. We use the LeNet-like convolutional

and deconvolutional autoencoders. On MNIST, the encoder contains two convolutional modules; $8 \times (5, 1, 2)$ convolutional layers and $4 \times (5, 1, 2)$ convolutional layers where a format is of (kernel size, stride, padding). The convolutional modules are followed by batch normalization, Leaky ReLU with $\alpha = 0.1$ and 2×2 max pooling. On CIFAR10, it contains three convolutional modules; $32 \times (5, 1, 2)$ convolutional layers followed by $64 \times (5, 1, 2)$ and $128 \times (5, 1, 2)$ convolutional layers. The dimension of the latent space is set to 32 for MNIST dataset, and 128 for CIFAR10 dataset. The decoders for both datasets contain symmetrical transposed convolutional modules to each encoder. We have used the reparameterization trick [6] with $L = 10$. Adam optimizer [15] is used for 1000 epochs with a learning rate of $5e-5$, a weight decay of $1e-4$ and batch size of 200. We did not conduct any data augmentations or pre-processing without normalizing the images to $[0, 1]$. The runs were performed with different β values from $\beta \in \{10, 5, 4, 3, 2, 1, 0.5, 0.4, 0.3, 0.2, 0.1\}$.

Results Fig. 3 shows the distortion-rate functions on datasets. Depending on β , we arrive at different R and D points. As expected, when β is increased, the resulting R values get decreased, which is what we have expected in Section 3. It is noted that as β decreased, $R + D$ gets increased meaning distortion-rate functions on these datasets are sub-optimal, which is due to the joint distribution $p(\mathbf{x}, \hat{\mathbf{x}})$. By means of using the powerful encoder-decoder or imposing an appropriate prior, we could shrink the gap to the optimal rate-distortion tradeoff more, but this is out of our scope. One can also find similar results from [8].

6.2 Anomaly Detection Performances

Baselines We compare our method, PGN, with other deep anomaly detection baselines for the anomaly detection task. We have considered three baselines using the reconstruction error as their anomaly score including naive autoencoder (AE), variational autoencoder (VAE) [6] and adversarial autoencoder (AAE) [28]. We also have considered Deep SVDD (DSVDD) [11] and GPND [35] as baselines.

Datasets and Experiment Protocol We have used MNIST, FMNIST, and CIFAR10 as benchmark datasets. We used official split training and test sets for all datasets. We took data of one class in the pre-split training set as our training set (thus, our training set consists of all normal instances) while our test set is the same as the pre-split test set. The test instances except for instances having the class label of the training set are deemed anomalies. MNIST, FMNIST, and CIFAR10 have 10 classes each. So we conducted 10 independent anomaly detection experiments for each dataset. The number of anomalies is about 9 times more than that of normal instances in test datasets.

Architectures and Experiment Settings The architectures of AE, AAE, VAE are the same as those in Section 6.1. For VAE, we also have used the reparameterization trick [6] with $L = 10$. AAEs train the model in adversarial learning and need the discriminator to discriminate between the generated latent variables of inputs and values sampled from the marginal prior. This discriminator consists of three fully-connected layers with the number of weights of $J \times 512$, 512×256 , and 256×1 , respectively, and each layer is followed by Leaky ReLU activations with $\alpha = 0.1$. Architectures of DSVDD and PGN are identical, the encoder of the autoencoder architecture. We also have considered *hypersphere collapse* as reported in [11] meaning that the model with bias terms can produce a trivial solution for DSVDD. We found that hypersphere collapse can also occur in our PGN setting, so we did not use any bias terms in both PGN and DSVDD. Also, PGN outputs the mean and variance values of isotropic multivariate Gaussian.

For PGN, we have used $T = 20$ for T stochastic forward passes with the dropout probability $p = 0.5$. It is noted that only mean outputs involve the MC dropout to consider the model uncertainty as described in Section 4 and Theorem 2.

It is reported that a pretraining with the autoencoder for DSVDD is helpful to enhance anomaly detection performance. But we did not conduct any pretraining for all baselines and PGN for fair comparisons. For AE, VAE, AAE, DSVDD and PGN, the Adam optimizer [15] is used with a weight decay of $1e-4$. The learning rate is initialized to $1e-4$ and reduced by a factor of 10 at 75th epochs for all datasets. We train for 100 epochs and compare under the AUROC values on test datasets. As stated in Section 2, we treat the anomalies are positive and normal instances are negative, so that the anomaly scores of the methods are used for calculating AUROC values without any fixed thresholds.

Table 1: Mean and std. dev. AUROCs [%] with 10 different seeds on MNIST (*Top*) and FMNIST (*Bottom*) datasets.

Normal class	GPND	DSVDD	AE	VAE	AAE	PGN (ours)
0	75.3±8.3	97.4±0.9	98.5±0.5	96.7±0.9	98.0±0.4	97.8±1.0
1	96.2±2.5	99.6±0.2	99.9±0.0	99.8±0.0	99.8±0.0	99.6±0.1
2	65.4±9.2	88.7±2.2	82.8±1.7	80.2±2.8	80.8±1.4	91.3±1.8
3	68.9±6.9	89.4±1.3	90.3±1.5	88.9±0.6	89.7±0.9	91.1±1.3
4	78.4±3.3	93.7±1.0	88.8±1.6	89.8±2.1	87.9±1.7	94.7±0.8
5	69.3±5.6	87.1±2.6	92.3±1.3	89.9±2.3	90.3±2.0	89.8±2.0
6	78.4±7.0	98.0±0.5	94.7±1.5	92.3±1.1	92.6±2.5	98.5±0.4
7	83.8±4.8	94.1±1.0	94.6±0.9	92.8±0.9	93.9±0.6	94.9±0.8
8	57.4±5.7	90.8±1.2	78.5±1.8	79.5±1.8	76.1±2.9	92.1±0.9
9	77.1±3.9	95.9±0.6	91.8±1.5	90.6±2.2	90.1±2.2	96.7±0.3
Avg.	75.0	93.5	91.2	90.1	89.9	94.7
T-shirt	77.2±7.9	90.4±1.1	88.1±0.5	87.5±0.6	87.9±0.6	90.4±3.4
Trouser	95.8±1.4	98.5±0.2	97.8±0.2	96.9±0.3	98.1±0.1	98.6±0.1
Pullover	78.1±5.6	85.8±3.1	83.7±0.6	83.8±0.8	80.7±1.6	86.1±5.2
Dress	85.5±4.4	92.4±1.5	90.8±0.5	89.3±0.7	90.5±0.5	93.1±1.3
Coat	77.8±4.3	89.2±1.4	86.7±0.5	85.1±0.8	86.6±0.6	87.7±6.5
Sandal	89.4±0.9	89.4±0.6	83.3±1.2	82.3±1.4	83.5±1.2	89.2±0.5
Shirt	76.3±3.5	80.6±1.7	78.7±0.3	79.0±0.6	77.5±0.4	80.3±2.3
Sneaker	95.3±1.5	98.6±0.1	97.6±0.1	96.9±0.1	97.6±0.1	98.7±0.1
Bag	68.1±4.2	91.1±1.9	75.2±1.4	75.3±2.1	75.5±2.3	92.2±5.5
Ankle boot	88.1±5.0	98.4±0.3	94.8±0.7	92.9±0.7	95.4±0.9	98.2±1.3
Avg.	83.2	91.4	87.7	86.9	87.3	91.5

For GPND, we have followed same experiment settings they recommended and provided⁵. For a data preprocessing, we only normalize the data to [0, 1] for both test and training datasets.

Table 2: Mean and std. dev. AUROCs [%] with 10 different seeds on CIFAR10 dataset.

Normal class	GPND	DSVDD	AE	VAE	AAE	PGN (ours)
Airplane	56.6±4.3	59.6±4.9	68.7±0.2	69.4±0.5	67.9±0.4	73.7±3.5
Automobile	55.2±3.8	56.4±2.4	40.0±0.3	44.6±0.4	41.6±0.3	55.1±2.8
Bird	55.1±2.4	64.9±1.8	65.0±0.2	65.1±0.2	65.1±0.2	64.3±3.0
Cat	56.5±3.3	53.4±0.7	55.8±0.3	53.8±0.2	54.2±0.2	56.7±2.5
Deer	69.2±1.8	72.4±3.1	67.2±0.1	68.5±0.2	68.1±0.3	70.5±2.1
Dog	53.9±2.0	53.5±2.4	56.2±0.2	53.9±0.2	54.4±0.2	60.3±3.8
Frog	70.0±5.4	74.2±2.7	55.4±0.3	61.1±0.5	58.7±0.4	72.1±3.5
Horse	58.4±2.1	54.3±2.4	44.9±0.4	46.2±0.3	45.1±0.2	54.9±2.9
Ship	64.8±4.0	67.6±2.2	74.5±0.3	74.5±0.4	73.4±0.3	75.5±2.8
Truck	54.5±4.7	60.9±2.8	41.8±0.2	45.2±0.5	42.7±0.3	60.8±2.2
Avg.	59.4	61.7	57.0	58.2	57.1	64.4

Comparison Results Table 1 shows the results on MNIST and FMNIST datasets. As shown, most of the cases, PGN gives the best results. AE gives competitive results among the reconstruction error based methods, AE, VAE, and AAE. AAE needs adversarial learning for training the generator (encoder) but it does not guarantee that it gives any meaningful differences with AE and/or VAE. VAEs incorporate the rate into the loss function but anomaly detection conducts only with the distortion (reconstruction error) so VAE also does not give any advanced results. Even though VAE gives the theoretical foundation on the autoencoder and can tweak the disentangled latent values to manipulate output easily, it does not have any advantages in anomaly detection tasks empirically. Table 2 gives the results on CIFAR10 dataset. It shows that PGN gives the most competitive results on this dataset as well. PGN and DSVDD empirically show very similar AUROC values along the

⁵<https://github.com/podgorskiy/GPND>

classes because they share the same theoretical formulation, but we can say PGN is slightly but meaningfully better than DSVDD. The main reason for this is that PGN detects anomalies with their involved model uncertainty metric, which can be captured by its anomaly score (Eq. 15).

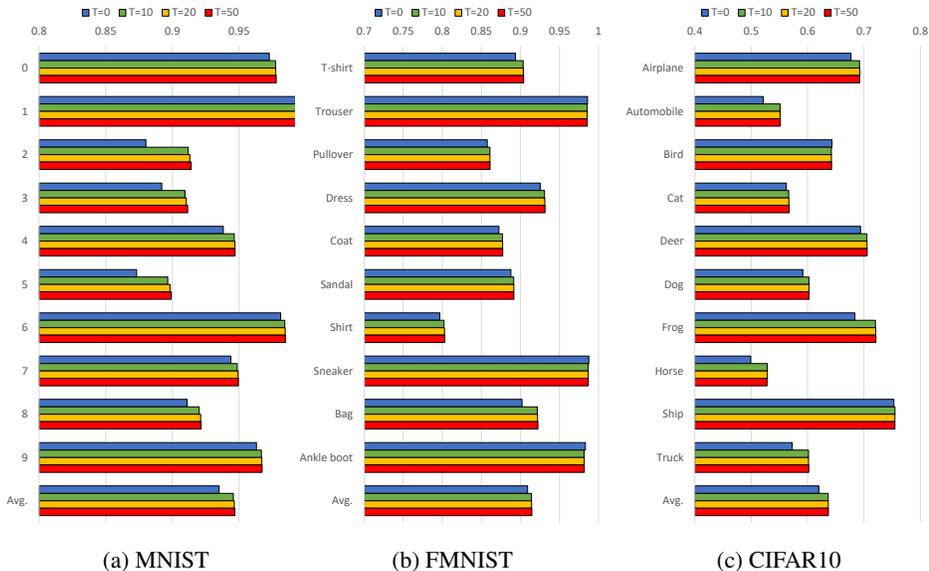


Figure 4: Ablation experiment results of model uncertainty metric of PGN on MNIST, FMNIST and CIFAR10 datasets. The values represent AUROC_s [%] with 10 different seeds.

Ablation Study To scrutinize the effectiveness of involving the model uncertainty we proposed, we performed an ablation study on three datasets. The results are shown in Fig. 4. In Fig. 4, $T = 0$ corresponds to the standard dropout meaning that we turn off stochastic forward passes when testing thus it does not consider the model uncertainty involved in the input instance. It shows that when $T \geq 10$ the mean and standard deviation values are very similar so we can say that using $T = 10$ is enough to leverage model uncertainty via MC dropout for these datasets. As a result, when we use MC dropout, the AUROC values on three datasets get increased with significant margins.

7 Conclusion

We would like to highlight that this work could connect the deep anomaly detection method with the theoretical foundations on variational inference and information theory. We propose PGN that can capture anomalies by means of estimating data distribution and shows better result by incorporating MC dropout based model uncertainty metric. We expect that this model could be better when equipped with more powerful architectures.

References

- [1] Simon Hawkins, Hongxing He, Graham Williams, and Rohan Baxter. Outlier detection using replicator neural networks. In *International Conference on Data Warehousing and Knowledge Discovery*, pages 170–180. Springer, 2002.
- [2] Mayu Sakurada and Takehisa Yairi. Anomaly detection using autoencoders with nonlinear dimensionality reduction. In *Proceedings of the MLSDA 2014 2nd Workshop on Machine Learning for Sensory Data Analysis*, pages 4–11, 2014.
- [3] Jinwon An and Sungzoon Cho. Variational autoencoder based anomaly detection using reconstruction probability. *Special Lecture on IE*, 2(1), 2015.
- [4] Erik Marchi, Fabio Vesperini, Florian Eyben, Stefano Squartini, and Björn Schuller. A novel approach for automatic acoustic novelty detection using a denoising autoencoder with bidirec-

- tional LSTM neural networks. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 1996–2000. IEEE, 2015.
- [5] Chong Zhou and Randy C Paffenroth. Anomaly detection with robust deep autoencoders. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 665–674, 2017.
- [6] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [7] Toby Berger. Rate-distortion theory. *Wiley Encyclopedia of Telecommunications*, 2003.
- [8] Alexander A Alemi, Ben Poole, Ian Fischer, Joshua V Dillon, Rif A Saurous, and Kevin Murphy. Fixing a broken ELBO. *arXiv preprint arXiv:1711.00464*, 2017.
- [9] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. β -VAE: Learning basic visual concepts with a constrained variational framework. *Iclr*, 2(5):6, 2017.
- [10] Bernhard Schölkopf, John C Platt, John Shawe-Taylor, Alex J Smola, and Robert C Williamson. Estimating the support of a high-dimensional distribution. *Neural computation*, 13(7):1443–1471, 2001.
- [11] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In *International conference on machine learning*, pages 4393–4402, 2018.
- [12] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in β -VAE. *arXiv preprint arXiv:1804.03599*, 2018.
- [13] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. *arXiv preprint arXiv:1802.05983*, 2018.
- [14] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2015.
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [16] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(Jul):2121–2159, 2011.
- [17] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311, 2018.
- [18] Seonho Park, Seung Hyun Jung, and Panos M Pardalos. Combining stochastic adaptive cubic regularization with negative curvature for nonconvex optimization. *Journal of Optimization Theory and Applications*, 184(3):953–971, 2020.
- [19] Matthew D Hoffman and Matthew J Johnson. ELBO surgery: yet another way to carve up the variational evidence lower bound. In *Workshop in Advances in Approximate Bayesian Inference, NIPS*, volume 1, page 2, 2016.
- [20] Jakub M Tomczak and Max Welling. VAE with a VampPrior. *arXiv preprint arXiv:1705.07120*, 2017.
- [21] Yarín Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. *arXiv preprint arXiv:1506.02142*, 2015.
- [22] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.

- [23] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- [24] Robert A Becker. The variance drain and Jensen’s inequality. *CAEPR Working Paper No. 2012-004*, 2012.
- [25] JG Liao and Arthur Berg. Sharpening Jensen’s inequality. *The American Statistician*, 73(3): 278–281, 2019.
- [26] David MJ Tax and Robert PW Duin. Support vector data description. *Machine learning*, 54(1): 45–66, 2004.
- [27] Hoang Anh Dau, Vic Ciesielski, and Andy Song. Anomaly detection using replicator neural networks trained on examples of one class. In *Asia-Pacific Conference on Simulated Evolution and Learning*, pages 311–322. Springer, 2014.
- [28] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.
- [29] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*, 11(Dec):3371–3408, 2010.
- [30] Jonathan Masci, Ueli Meier, Dan Cireşan, and Jürgen Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. In *International conference on artificial neural networks*, pages 52–59. Springer, 2011.
- [31] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [32] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Ursula Schmidt-Erfurth, and Georg Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *International conference on information processing in medical imaging*, pages 146–157. Springer, 2017.
- [33] Pramuditha Perera, Ramesh Nallapati, and Bing Xiang. OCGAN: One-class novelty detection using GANs with constrained latent representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2898–2906, 2019.
- [34] Raghavendra Chalapathy, Aditya Krishna Menon, and Sanjay Chawla. Anomaly detection using one-class neural networks. *arXiv preprint arXiv:1802.06360*, 2018.
- [35] Stanislav Pidhorskyi, Ranya Almoheisen, and Gianfranco Doretto. Generative probabilistic novelty detection with adversarial autoencoders. In *Advances in neural information processing systems*, pages 6822–6833, 2018.
- [36] Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000.
- [37] Rob Brekelmans, Daniel Moyer, Aram Galstyan, and Greg Ver Steeg. Exact rate-distortion in autoencoders via echo noise. In *Advances in Neural Information Processing Systems*, pages 3884–3895, 2019.
- [38] Yochai Blau and Tomer Michaeli. Rethinking lossy compression: The rate-distortion-perception tradeoff. *arXiv preprint arXiv:1901.07821*, 2019.
- [39] Luis A Lastras. Information theoretic lower bounds on negative log likelihood. *arXiv preprint arXiv:1904.06395*, 2019.
- [40] Yann LeCun. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [41] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

[42] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. The CIFAR-10 dataset. *online: <http://www.cs.toronto.edu/kriz/cifar.html>*, 55, 2014.