

Anomaly Detection Based on Unsupervised Disentangled Representation Learning in Combination with Manifold Learning

Xiaoyan Li

*Electrical Engineering & Computer Science
University of Ottawa
Ottawa, Canada
xli343@uottawa.ca*

Iluju Kiringa

*Electrical Engineering & Computer Science
University of Ottawa
Ottawa, Canada
Iluju.Kiringa@uottawa.ca*

Tet Yeap

*Electrical Engineering & Computer Science
University of Ottawa
Ottawa, Canada
tyeap@uottawa.ca*

Xiaodan Zhu

*Electrical & Computer Engineering
Queen's University
Kingston, Canada
xiaodan.zhu@queensu.ca*

Yifeng Li

*Computer Science
Brock University
St. Catharines, Canada
yli2@brocku.ca*

Abstract—Identifying anomalous samples from highly complex and unstructured data is a crucial but challenging task in a variety of intelligent systems. In this paper, we present a novel deep anomaly detection framework named AnoDM (standing for Anomaly detection based on unsupervised Disentangled representation learning and Manifold learning). The disentanglement learning is currently implemented by β -VAE for automatically discovering interpretable factorized latent representations in a completely unsupervised manner. The manifold learning is realized by t-SNE for projecting the latent representations to a 2D map. We define a new anomaly score function by combining β -VAE's reconstruction error in the raw feature space and local density estimation in the t-SNE space. AnoDM was evaluated on both image and time-series data and achieved better results than models that use just one of the two measures and other deep learning methods.

Index Terms—Anomaly detection, disentangled representation learning, manifold learning.

I. INTRODUCTION

Detecting anomalies in data flow of modern intelligent systems is an important but challenging problem. Formally speaking, anomaly detection problems can be statistically viewed as identifying outliers having low probabilities from the modelling of data distribution $p(\mathbf{x})$. Practically, since statistical modelling of the data is often difficult, it degenerates to domain description [1] or supervised prediction [2] problems in some cases. The exact explanation of an anomalous data point depends on the specific domain of focus. In data centers, it probably indicates an attempt of cyber intrusion. In recognition systems, it could be an adversarial attack. In biomedical information systems, it means possible onset of certain diseases. In Internet of Things (IoT) systems, it may represent a hardware failure or alarming event captured by

sensors. An anomalous sample is not always associated with negativity. Sometimes, it leads to novel discoveries in scientific explorations.

However, from the data analytics perspective, anomaly detection is a difficult task due to the following reasons. (1) Many forms of data, e.g., images, text, and other types of sequences, are often highly unstructured and complex. How can these data be well represented and high-level information be extracted by an algorithm? (2) The sample sizes of modern data sets are often extremely large and most of them are unlabelled. Unfortunately, traditional methods do not scale and perform well on these data. (3) When data of multiple modalities are naturally available for same events in a system, a robust and precise algorithm needs to be designed to integrate these information for system diagnosis or decision making. (4) Many intelligent systems, such as IoTs, require real-time detection and reaction of abnormal events to avoid costly and irrevocable damages. Thus, anomaly monitoring algorithms to be designed in these platforms must be highly efficient. In summary, anomaly detection raises challenges in representability, scalability, multimodality, and time complexity.

Deep learning [3] offers great potentials to overcome these challenges. (1) Representation learning mechanisms (such as convolution for images, embedding for discrete symbols, and recurrence for time-series) have been developed in supervised and unsupervised deep models to consider the nature of specific types of input samples and encode them into vectors of continuous values as corresponding latent representations. (2) Most deep learning models are trained using stochastic gradient descent that splits a giant training set into mini-batches. Thus, learning becomes unrestricted and blessed by a large sample size. Particularly, stochastic variational inference [4], [5] has successfully enabled scalable learning and infer-

ence for deep generative models (DGMs) on a vast amount of unlabelled data. (3) The development of deep learning programming packages, such as PyTorch [6] and TensorFlow [7], greatly eases the assembly of multiple network components (corresponding to different modalities) together for multimodal representation learning [8]. (4) Once a deep model is learned, the inference or encoding step is very efficient, thanks to the highly parallel computing architectures and techniques.

In some applications, if the domain of anomalous and normal samples is well defined, anomaly detection can be reduced to binary classification problems. However, in many situations, either the domain of anomalous samples cannot be fully understood or modelled, or the domain of the normal samples is too complicated to be modelled in one class. DGMs are more suitable than supervised methods in such cases. DGMs are concerned with the joint distribution of visible and latent variables with a hierarchy of stochastic (and deterministic) layers. With proper emphasis on disentanglement of latent representations, DGMs have the potential of dissecting hidden factors that are key to sample generation. Unsupervised disentangled representation learning [9] renders several benefits. (1) It helps better understand our data, providing a path towards explainable AI. (2) It gives a better control on the generation process of novel samples. (3) The disentanglement of latent factors may provide an opportunity to distinguish anomalies based on the landscape of latent space, which is our interest in this paper. It has been shown that the likelihood of a data point $p(\mathbf{x})$ estimated in DGM is not a reliable measure for detecting abnormal samples [10]. Instead, reconstruction error is widely used as an anomaly score function [11].

As a variant of variational autoencoder (VAE) [12], β -VAE [13] is designed for unsupervised discovery of interpretable factorized latent representations from raw image data. An adjustable hyperparameter β is introduced to balance the extent of learning constraints (a limit on the capacity of the latent information channel and an emphasis on learning statistically independent latent factors) and reconstruction accuracy. It was demonstrated that β -VAE with appropriately tuned value of β (when $\beta > 1$) qualitatively outperforms VAE (when $\beta = 1$, β -VAE is exactly VAE). [14] proposed a modification to the training regime of β -VAE by progressively increasing the information capacity of the latent code during training. This modification facilitates the robust learning of disentangled representations in β -VAE, without the previous trade-off in the reconstruction accuracy. [15] introduced a reformulation of β -VAE for $0 < \beta < 1$. They argued that, within in this range, training β -VAE is equivalent to optimizing an approximate log-marginal likelihood bound of VAE under an implicit prior.

Manifold learning is a family of nonlinear dimensionality reduction techniques. The t-distributed stochastic neighbor embedding (t-SNE) [16] is an unsupervised manifold learning method primarily used for data exploration and visualization by approximating high-dimensional data distribution using a two or three-dimensional map that could preserve local and certain global structures of the data. The use of t-SNE for anomaly detection has been sceptical [16]. However, no com-

prehensive investigation has been made in this topic. Taking advantages of both disentangled representation learning (using β -VAE as an implementation) and low-dimensional manifold learning (using t-SNE as an implementation), we propose a novel anomaly detection approach named **AnoDM**, standing for *Anomaly detection based on unsupervised Disentangled representation learning and Manifold learning*. We introduce a new anomaly score function by combining: (1) β -VAE’s reconstruction error, and (2) distances between latent representations of test points and training points in t-SNE map. AnoDM is a general framework, thus any disentangled representation learning and manifold learning techniques can be applied. The choice of a lower-level encoding scheme in β -VAE depends on data type of interest. For image data, deterministic convolutional network (CNN) is used in the encoder. In case of time series (sequence) data, we design an improved version of β -VAE by replacing CNN with temporal convolutional network (TCN) [17], a generic architecture for convolutional sequence prediction, in the encoder. We incorporate TCN as part of the encoder, because [17] have shown that TCN outperforms canonical recurrent networks such as LSTMs [18] across a range of supervised learning tasks and recommended that CNN should be regarded as the first method to try for sequence modeling tasks. Regarding the decoding architecture, we simply choose CNN, because by choosing a simpler CNN architecture as a part of the decoder, the model can achieve a comparable even better performance but take much less running time.

The contributions of this paper are summarized as follows. (1) We comprehensively explore the capacity of unsupervised disentangled representation learning, using β -VAE as an implementation, for anomaly detection. (2) We thoroughly investigate the potential of manifold learning for outlier identification by taking the disentangled latent representations from β -VAE as input to t-SNE. To the best of our knowledge, this is the first attempt to explore t-SNE for anomaly detection. (3) For sequence anomaly detection, instead of using prevailed recurrent networks (such as LSTM), as a practical contribution, we adopt an improved convolution architecture (TCN) to capture the temporal dependency in the encoder in unsupervised way.

II. RELATED WORK

In the big data era, the development of deep learning models, especially DGMs, flourishes due to the need of modelling and analyzing massive amount of unstructured data (such as images, time-series, graphs, text, etc.) generated in many application domains. Designing DGM-based solutions for anomaly detection becomes an important topic. Since DGMs, such as VAE [12], and deep belief net (DBN) [19], [20], aim at modelling the joint distribution of visible and latent variables (that is $p(\mathbf{x}, \mathbf{h})$), their likelihood $p(\mathbf{x})$ by marginalizing out \mathbf{h} may serve as an abnormality indicator. However, unlike exponential family restricted Boltzmann machines (exp-RBMs) [21], exact likelihood is unavailable for most DGMs. Alternatively, reconstruction error serves as an abnormality measure based on the

intuition that out-of-distribution samples can be reconstructed badly [11]. Some deep hybrid methods (e.g. VAE+OCSVM [22] and DBN+OCSVM [23]), successfully combine classical one-class support vector machine (OCSVM; or kernel-based support vector domain description (SVDD)) with DGMs by using DGMs to learn latent representations of samples and using OCSVM to detect abnormal data points. However, these methods face the challenge of scalability, because the size of kernel matrices in dual form of SVDD is quadratic of sample size.

Generative adversarial net (GAN) has also been applied for anomaly detection [24]. Since there is no encoder in GAN, [25] presented the ADGAN algorithm based on the availability of a good representation of a sample in latent space of its generator by assuming that the generator is able to effectively capture the distribution of the training data. [26] proposed the GAN-AD method for cyber-physical systems (CPSs). It distinguishes fake data from actual data by taking into consideration of both discrimination loss calculated by the trained discriminator and residual loss between reconstructed and actual test data.

Furthermore, DGM-based algorithms are also devised to detect anomaly problem on sequence data (e.g. LSTM-VAE [27] and GAN-AD [26]). Conventionally, canonical recurrent networks (such as LSTM and GRUs [28]) are considered as the dedicated methods for sequence modeling. Some recent studies have also claimed that there was no architecture that could consistently beat LSTM in some typical sequence modelling tasks [29]–[31]. On the other hand, some other researchers insist that CNN [32] should be considered as more appropriate choice for sequences. Inspired by more recent CNN-based sequence modelling (such as machine translation [33], [34] and language modeling [35]), [17] conducted a systematic evaluation of generic convolutional and recurrent architectures for sequence modelling across a broad range of tasks that are commonly used to benchmark recurrent networks, and concluded that convolutional networks, rather than recurrent networks, should be respected as a “natural starting point for sequence modelling tasks”.

III. METHOD

In this paper, we propose a novel generic anomaly detection framework named AnoDM, which the first time combines unsupervised disentangled representation learning (implemented using β -VAE as an example) and low-dimensional manifold learning (currently using t-SNE as implementation) together to detect outliers via effectively taking the advantages of reconstruction at raw feature space and disentangled latent distribution in t-SNE map. Fig. 1 shows the architecture of AnoDM which includes two main phases: (1) unsupervised disentangled representation learning and (2) anomaly detector. After β -VAE is learned using unlabelled training normal samples, it then can be employed by the anomaly detector to efficiently obtain latent encoding and reconstructed version of a sample. Once latent embeddings of both training samples (or a representative subset from the training data) and a test sample

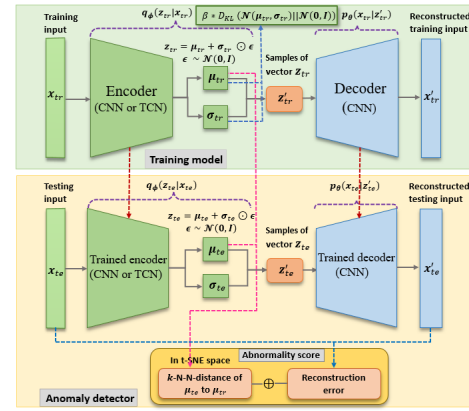


Fig. 1: Architecture of AnoDM implemented by β -VAE and t-SNE for anomaly detection. First, β -VAE is learned using normal training data (upper part of the framework). Then it is employed by the anomaly detector (lower part of the framework) to efficiently obtain latent encodings of training samples (or a representative subset from the training set) and test samples, and corresponding reconstructed versions using the decoder. Meanwhile, t-SNE is used to map the deterministic latent embeddings (μ_{tr} for training and μ_{te} for test samples) to the 2D space (we call it t-SNE space or map), such that the average distance between the 2D representation of a test sample and its k nearest neighbors from the 2D representations of training samples is calculated. Finally, the distance is combined with the reconstruction error of the test sample to define its anomaly score.

are obtained, t-SNE is used to map the latent representations of these samples further to the 2-dimensional space (called t-SNE space or map), such that the average distance between the 2D representation of the test sample and its k nearest neighbors from the 2D representations of training samples is calculated. Finally, this distance is combined with the reconstruction error of the test sample to define its anomaly score. The essential parts of this framework are discussed below in details. Full AnoDM approach is given in Algorithm 1.

A. Unsupervised Disentangled Representation Learning

The unsupervised disentangled representation learning component in our architecture is implemented but not limited by β -VAE [13]–[15], [36]. The objective function to be maximized for β -VAE is defined as [14]:

$$\mathcal{L} = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \beta |D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) - C|. \quad (1)$$

The first term of this objective function corresponds to reconstruction error in the raw feature space. The KL divergence characterizes the discrepancy between approximate posterior and isotropic prior of latent representations. A small discrepancy between them indicates high disentanglement of latent representations in independent variables. C is a hyperparameter which is used to improve the quality of reconstructed images. The loss function of the original β -VAE proposed in [13] does not have this hyperparameter. The value of β trades off reconstruction error and disentanglement. Unlike [13] and [14], we consider $\beta > 0$ rather than just $\beta > 1$, because it is unnecessary to bound the value of β by 1, $\beta > 0$ allows us search for a more appropriate disentanglement. The special case $\beta = 0$ could make the model learning

Algorithm 1: AnoDM Algorithm

Result: Anomaly scores of test samples

Inputs: \mathbf{X}_{tr} : training samples, \mathbf{X}_{te} : test samples,
 $\beta > 0$: hyperparameter for β -VAE

```
1 while epoch no more than training iterations do
2   Encoder net maps  $\mathbf{X}_{\text{tr}}$  into  $\boldsymbol{\mu}_{\text{tr}}$  and  $\boldsymbol{\sigma}_{\text{tr}}$ ;
3    $\mathbf{Z}_{\text{tr}} = \boldsymbol{\mu}_{\text{tr}} + \boldsymbol{\sigma}_{\text{tr}} \odot \boldsymbol{\epsilon}$ ,  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathcal{I})$ ;
4   Decoder net reconstructs  $\mathbf{X}_{\text{tr}}$  to  $\mathbf{X}'_{\text{tr}}$  using  $\mathbf{Z}_{\text{tr}}$ ;
5   Update  $\beta$ -VAE's parameters  $\boldsymbol{\theta}$  and  $\boldsymbol{\phi}$ ;
6 end
7 For  $\mathbf{X}_{\text{tr}}$  and  $\mathbf{X}_{\text{te}}$ , obtain  $\boldsymbol{\mu}_{\text{tr}}$  and  $\boldsymbol{\mu}_{\text{te}}$  respectively using
  trained  $\beta$ -VAE;
8 Use t-SNE to map  $\boldsymbol{\mu}_{\text{tr}}$  and  $\boldsymbol{\mu}_{\text{te}}$  to 2D representations
   $\mathbf{l}_{\text{tr}}$  and  $\mathbf{l}_{\text{te}}$ ;
9 for  $\mathbf{x}_{\text{te}}^{(i)}$  within  $\mathbf{X}_{\text{te}}$  do
10   $\mathcal{D}_{\text{RE}}(\mathbf{x}_{\text{te}}^{(i)}) \triangleq \text{NSE}(\mathbf{x}_{\text{te}}^{(i)}, \mathbf{x}'_{\text{te}}{}^{(i)}) = \frac{\|\mathbf{x}_{\text{te}}^{(i)} - \mathbf{x}'_{\text{te}}{}^{(i)}\|_2^2}{\|\mathbf{x}_{\text{te}}^{(i)}\|_2}$ ;
    // reconstruction error
11   $\mathcal{D}_{\text{tSNE}}^k(\mathbf{x}_{\text{te}}^{(i)}) \triangleq \frac{1}{k} \sum_{j \in N(i,k)} \|\mathbf{l}_{\text{te}}^{(i)} - \mathbf{l}_{\text{tr}}^{(j)}\|_2$ ;
    //  $N(i,k)$  is the set of indices of  $\mathbf{l}_{\text{te}}^{(i)}$ 's  $k$ 
    nearest neighbors from  $\mathbf{l}_{\text{tr}}$ 
12   $\mathcal{S}_{\beta\text{VAE+tSNE}}(\mathbf{x}_{\text{te}}^{(i)}) = \alpha \mathcal{D}_{\text{RE}} + (1 - \alpha) \mathcal{D}_{\text{tSNE}}^k$ ;
    //  $\alpha \in [0, 1]$ 
13 end
```

very unstable, because the variance of inference distribution loses control. It worth highlighting that other unsupervised disentanglement models can be used as well in AnoDM. For example, [36] interpreted disentanglement as decomposition instead of independence by adding an additional regularization term to reduce the discrepancy between the aggregate posterior and a desired structured prior. However, designing a properly structured prior could be practically challenging. The adoption of β -VAE in our framework is sufficient to prove the concept that unsupervised disentanglement helps anomaly detection.

B. Effectivity of t-SNE Algorithm

In addition to β -VAE's reconstruction error, we use the average distance between a test sample and its k -nearest neighbors from the collection of training samples in the t-SNE [16] map to score the outlieriness of a test sample. As t-SNE is significantly influenced by perplexity, the nature of complexity in data distributions makes it impossible to utilize a uniform criteria to define optimal perplexity for all data. Moreover, [37] mentioned several weaknesses of t-SNE, for examples, (1) it naturally expands dense clusters and contracts sparse ones, evening out cluster sizes, and (2) distances between clusters might not reflect global geometry. However, it is likely that k -nearest neighbors still work for local clumps, because, with proper value of perplexity, local topological information of latent distributions can be preserved by the t-SNE plot. Thus, in t-SNE space, the measure of k -nearest neighbors is more suitable than full density estimation which is very sensitive with the sizes of clusters. Furthermore, as to be shown in Section IV, distancing in the 2D t-SNE map could be more

robust than distancing in β -VAE's latent space where many non-determinative factors may influence the calculation of distances. Finally, it worth clarifying that we do not directly learn a 2D latent representation from β -VAE, because it will bottleneck too much the information flow for reconstruction. Instead, a lower-dimensional representation is learned through t-SNE for density estimation only.

C. Deterministic or Stochastic Latent Representations for t-SNE

Either the mean $\boldsymbol{\mu}$ or a sample \mathbf{z} from the approximate inference distribution $q(\mathbf{z}|\mathbf{x})$ can be passed to t-SNE to calculate the k -NN distance of a test sample. There is trivial difference between performances achieved by these two methods in our framework. Generally, the $\boldsymbol{\mu}$ -based method achieved slight better performance. The comparison of these two methods can be found in Table I. Furthermore, from the latent representations' t-SNE maps (see Fig. 2), one can interestingly see that when β is small (not overly large), the t-SNE maps for both methods are quite similar. As β becomes overly large, some normal classes can still form their own clusters (even though some similar classes, such as classes 3 and 8 in MNIST, tend to mingle together) in $\boldsymbol{\mu}$ -based method, but in \mathbf{z} -based method all classes are entangled with each other. Same phenomena can be observed on the other datasets. Therefore, the $\boldsymbol{\mu}$ -based method is used in current design of AnoDM.

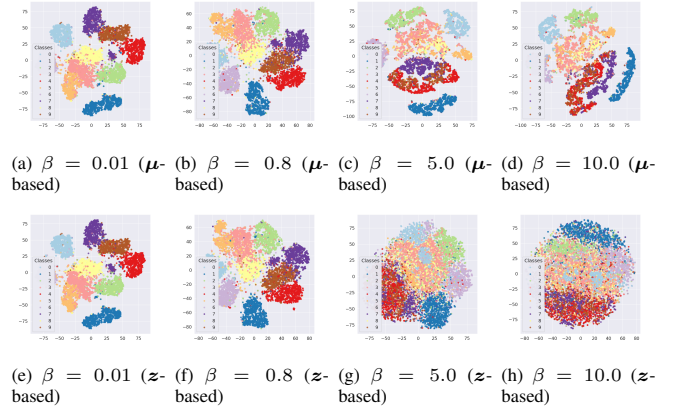


Fig. 2: Comparison of $\boldsymbol{\mu}$ -based method or \mathbf{z} -based method on MNIST data (anomalous class is 3) for inferring latent representations that are visualized in t-SNE map.

D. TCN Encoder for Unsupervised Sequence Modelling

[17] distilled superior design in convolutional network into a simple architecture and referred it as a temporal convolutional network (TCN) with two distinctive characteristics: (1) the convolutions in the architecture are causal, and (2) the architecture can take a sequence of any length and map it to an output vector of fixed length, just as with an RNN. [17] also explained that TCNs capture significantly longer history than recurrent networks. Inspired by [17], we replace CNN with TCN in the encoder of β -VAE when evaluating the proposed

TABLE I: AuROCs for both μ -based and z -based techniques.

Dataset	Class	μ -based			z -based		
		auROC	β	α	auROC	β	α
MNIST	0	0.984	0.01	0.8	0.985	0.8	0.1
	1	0.986	0.01	0.6	0.987	0.01	0.6
	2	0.990	0.4	0.9	0.991	0.2	0.9
	3	0.969	0.05	0.9	0.968	0.05	0.9
	4	0.974	0.01	0.95	0.975	0.1	0.95
	5	0.975	0.01	0.95	0.976	0.01	0.95
	6	0.983	0.01	0.8	0.980	0.01	0.8
	7	0.975	0.01	0.9	0.977	0.01	0.9
	8	0.980	0.01	0.9	0.982	0.01	0.9
	9	0.928	0.4	0.8	0.925	0.05	0.8
avg.		0.974	-	-	0.975	-	-
Fashion-MNIST	0	0.844	0.05	0.2	0.840	0.05	0.0
	1	0.977	0.01	0.8	0.978	0.01	0.8
	2	0.783	0.01	0.0	0.783	0.05	0.0
	3	0.886	0.01	0.8	0.884	0.01	0.8
	4	0.760	0.1	0.0	0.763	0.3	0.0
	5	0.990	0.2	0.95	0.990	0.2	0.95
	6	0.713	0.05	0.0	0.709	0.05	0.0
	7	0.952	0.01	0.8	0.950	0.01	0.8
	8	0.980	0.05	0.8	0.980	0.05	0.8
	9	0.940	0.3	0.8	0.944	0.3	0.7
avg.		0.883	-	-	0.882	-	-
CIFAR-10	0	0.635	0.4	0.0	0.634	0.2	0.0
	1	0.752	0.6	0.99	0.754	0.05	0.99
	2	0.589	1.0	0.0	0.558	0.7	0.0
	3	0.608	10.0	0.95	0.606	10.0	0.99
	4	0.564	0.01	0.0	0.563	0.4	0.0
	5	0.638	0.4	0.95	0.627	0.7	0.95
	6	0.600	0.3	0.0	0.590	0.7	0.0
	7	0.648	0.1	0.95	0.644	0.01	0.95
	8	0.642	0.3	0.0	0.624	1.2	0.0
	9	0.717	1.0	0.95	0.718	0.01	0.95
avg.		0.639	-	-	0.632	-	-
Small-Norb	0	0.512	0.1	0.0	0.520	7.0	0.0
	1	0.656	0.01	0.0	0.647	0.01	0.0
	2	0.771	0.05	1.0	0.771	0.05	1.0
	3	0.581	10.0	1.0	0.581	10.0	1.0
	4	0.564	0.5	1.0	0.564	0.5	1.0
avg.		0.617	-	-	0.617	-	-
ECG (Arrhythmia)	0	0.911	0.05	0.95	0.906	0.01	0.95
	1	0.924	0.01	0.95	0.925	0.01	0.95
	2	0.970	0.01	0.99	0.970	0.01	0.99
	3	0.905	0.01	0.95	0.910	0.01	0.95
	4	0.991	0.01	0.99	0.991	0.01	0.99
avg.		0.940	-	-	0.940	-	-

AnoDM framework on time-series data, while we still use CNN in the decoder, because our preliminary experiments demonstrated that keeping decoder as simpler CNN can help achieve comparable even better results, and take much less computing time. In the architecture of TCN, the kernel size is set to 4 and dilation factors are set to [1, 2, 4, 8, 16, 32]. In Section IV, the comparison among TCN, CNN, and LSTM encoders in our framework also shows that TCN outperforms CNN and particularly LSTM to a great extent for ECG signal anomaly detection.

E. Anomaly Score Function in AnoDM

In the anomaly detector, the reconstruction error of a test sample in the original feature space and the average distance from its k -nearest-neighbors in training samples within the 2D t-SNE map are combined as a final anomaly score function:

$$\mathcal{S}_{\beta\text{-VAE+tSNE}}(\mathbf{x}_{\text{te}}) = \alpha \mathcal{D}_{\text{RE}}(\mathbf{x}_{\text{te}}) + (1 - \alpha) \mathcal{D}_{\text{tSNE}}^k(\mathbf{x}_{\text{te}}), \quad (2)$$

where the first term is defined using normalized squared error (NSE):

$$\mathcal{D}_{\text{RE}}(\mathbf{x}_{\text{te}}) \triangleq \text{NSE}(\mathbf{x}_{\text{te}}, \mathbf{x}'_{\text{te}}) = \frac{\|\mathbf{x}_{\text{te}} - \mathbf{x}'_{\text{te}}\|_2^2}{\|\mathbf{x}_{\text{te}}\|_2}, \quad (3)$$

where \mathbf{x}_{te} is a test sample, and \mathbf{x}'_{te} is its reconstructed version by sending the stochastic latent encoding through the decoder

of β -VAE. The second term in Equation (2) is defined using β -VAE's deterministic latent encoding (mean from the encoder of learned β -VAE) as input to t-SNE:

$$\mathcal{D}_{\text{tSNE}}^k(\mathbf{x}_{\text{te}}^{(i)}) \triangleq \frac{1}{k} \sum_{j \in N(i,k)} \|\mathbf{l}_{\text{te}}^{(i)} - \mathbf{l}_{\text{tr}}^{(j)}\|_2, \quad (4)$$

where $\mathbf{l}_{\text{te}}^{(i)}$ is the 2D representation of the i -th test sample in t-SNE map, $N(i, k)$ is the set of indices of $\mathbf{l}_{\text{te}}^{(i)}$'s k nearest neighbors from training samples' 2D representations \mathbf{l}_{tr} in t-SNE map. In Equation (2), $\alpha \in [0, 1]$ is the combination hyperparameter such that the two terms can effectively complement each other. To allow the anomaly score function to achieve its full potential, α value should be sensitively searched, because the values of \mathcal{D}_{RE} and $\mathcal{D}_{\text{tSNE}}^k$ can stay at different magnitudes, a very small change of α value may dramatically alter the contributions of these two terms. Alternatively, the distance score in Equation 4 can be normalized by average distance of training samples in t-SNE map, which may alleviate the magnitude difference, thus ease search of optimal value of α . The use of this normalized distance is investigated in Section IV-C.

IV. EXPERIMENTS

We evaluated our framework on four public image datasets, including MNIST [38], Fashion-MNIST [39], Small-Norb [40] and CIFAR-10 [41], as well as one collection of ECG heartbeat categorization time-series data named Arrhythmia [42]. The detail of β -VAE architecture is given in Table III in appendix. The number of epochs was set to 20 for experiments on MNIST and Fashion-MNIST datasets, and 50 for CIFAR-10, Small-Norb, and Arrhythmia datasets; batch size was set to 100 for experiments on all these datasets. When using t-SNE, the dimension of t-SNE map was set as 2 for all datasets, perplexity 30, the learning rate 200, and maximum number of iterations 1000. The value of α in the anomaly score function is searched from set $\{0.0, 0.005, 0.01, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.95, 0.99, 1.0\}$. We used $k = 1$ for calculating k -nearest-neighbors distance in t-SNE map (we also tried to set k to 3 or 5, the results were similar).

A. Comparison with CapsNet, GANs, and VAE

We compared AnoDM with state-of-the-art algorithms including a supervised method – CapsNet [43], and two types of generative models – GANs (including AnoGAN and ADGAN) [25] and β -VAE (implemented by setting $\alpha = 1$ thus using only reconstruction error as anomaly score). As shown in Table II, on average, AnoDM achieved either comparable (on MNIST) or better (on Fashion-MNIST and CIFAR-10) performance in terms of receiver operating characteristic curve (auROC). On Fashion-MNIST, CapsNet (prediction-probability-based), as the best benchmark method, obtained an average auROC of 0.765, while AnoDM achieved 0.883. On MNIST, both AnoDM and CapsNet obtained the highest performance. However, CapsNet is a supervised method that

takes advantages of class information, while ours is completely unsupervised that is more suitable in many practices as class information is either incomplete or unavailable. Furthermore, by comparing AnoDM with β -VAE that only considers reconstruction error as anomaly score, AnoDM dramatically improved the performance in all cases. In other words, t-SNE makes a prominent contribution to improve β -VAE for anomaly detection problems. However, all generative models did not work well on Small-Norb, mainly because these models used convolution to extract features from image, but convolution is only able to capture translation but not other affine transformations. Although CapsNet learns these transformations as a supervised method, it worth exploring unsupervised learning of affine transformations as a future topic.

TABLE II: Performance of AnoDM in comparison with other methods in terms of auROC on image data. The results for AnoDM were obtained by either the μ -based or z -based algorithm. In the column for β -VAE, only reconstruction error is used as anomaly score. The results for AnoGAN and ADGAN were obtained from [25]. The results for CapsNet were obtained from [43].

Dataset	Class	CapsNet		AnoGAN	ADGAN	β -VAE	AnoDM
		PP-based	RE-based				
MNIST	0	0.998	0.947	0.990	0.999	0.890	0.985
	1	0.990	0.907	0.998	0.992	0.841	0.987
	2	0.984	0.970	0.888	0.968	0.967	0.991
	3	0.976	0.949	0.913	0.953	0.947	0.969
	4	0.935	0.872	0.944	0.960	0.968	0.975
	5	0.970	0.966	0.912	0.955	0.966	0.976
	6	0.942	0.909	0.925	0.980	0.907	0.983
	7	0.987	0.934	0.964	0.950	0.899	0.977
	8	0.993	0.929	0.883	0.959	0.946	0.982
	9	0.990	0.871	0.958	0.965	0.794	0.928
	avg.	0.977	0.925	0.937	0.968	0.913	0.975
Fashion-MNIST	0	0.620	0.454	–	–	0.500	0.844
	1	0.851	0.871	–	–	0.860	0.978
	2	0.818	0.486	–	–	0.459	0.783
	3	0.895	0.693	–	–	0.730	0.886
	4	0.790	0.394	–	–	0.379	0.763
	5	0.691	0.982	–	–	0.985	0.990
	6	0.801	0.480	–	–	0.501	0.713
	7	0.619	0.787	–	–	0.842	0.952
	8	0.912	0.885	–	–	0.876	0.980
	9	0.656	0.754	–	–	0.701	0.944
	avg.	0.765	0.679	–	–	0.683	0.883
CIFAR-10	0	0.622	0.371	0.610	0.661	0.368	0.635
	1	0.455	0.737	0.565	0.435	0.746	0.754
	2	0.671	0.421	0.648	0.636	0.397	0.589
	3	0.675	0.588	0.528	0.488	0.604	0.608
	4	0.683	0.388	0.670	0.794	0.387	0.564
	5	0.635	0.601	0.592	0.640	0.611	0.638
	6	0.727	0.491	0.625	0.685	0.500	0.600
	7	0.673	0.631	0.576	0.559	0.614	0.648
	8	0.513	0.410	0.723	0.798	0.399	0.642
	9	0.466	0.671	0.582	0.643	0.698	0.718
	avg.	0.612	0.531	0.612	0.634	0.532	0.640

B. Impact of Beta to Performance

As discussed in [13], β (> 1) functions as a controller to encourage most efficient latent representation learning via limiting the capacity of latent information channel. [36] however interpreted the objective of β as tuning a proper level of overlap of encodings by working with another term that regularizes the divergence of the aggregate posterior $q_\phi(z)$ and the desired prior $p(z)$. The main intuition is that purely increasing β induces too much overlap which actually discourages the disentanglement of data (information which is necessary for expressing desired structure is lost). [13] demonstrated that β -VAE with $\beta > 1$ leads to interesting results when learning

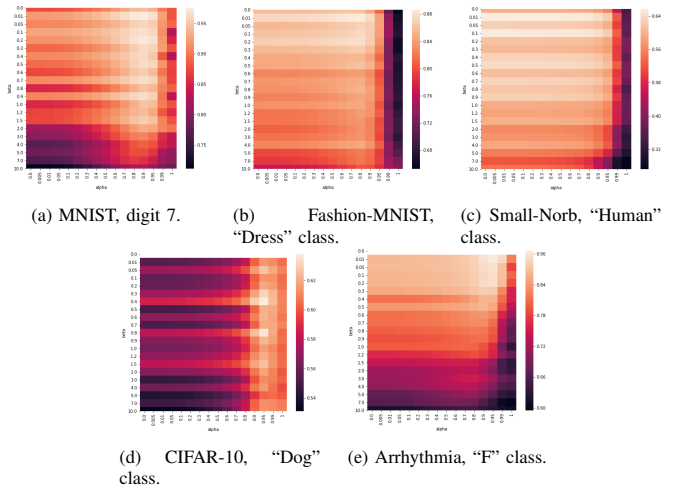


Fig. 3: Performances (measured in terms of auROC) of AnoDM evaluated on five datasets: MNIST, Fashion-MNIST, Small-Norb, CIFAR-10 and Arrhythmia. On each dataset, the anomalous class is indicated in the corresponding subcaption while treating the rest classes as normal classes. Note that $\beta = 0$ doesn't work for CIFAR10 and Arrhythmia, because it makes learning highly unstable. As displayed in (d) and (e), missing values were indicated in white color at the top of corresponding heatmaps for $\beta = 0$.

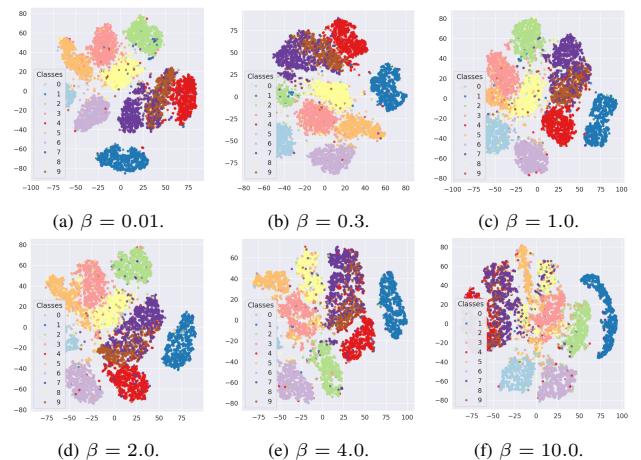


Fig. 4: The impact of β 's value to t-SNE map of latent representations of MNIST samples. Class 7 is treated as anomalous class data. Each map displays 10000 data points identical in all maps including 5000 training data points and 5000 test data points.

interpretable factorized latent representations on a variety of datasets. Surprisingly, our investigation demonstrates that by setting $0 < \beta < 1$, it actually achieved the state-of-the-art results for anomaly detection problems on a range of datasets, such as MNIST, Fashion-MNIST and Arrhythmia. Fig. 3 illustrates the impact of values of β and α in our anomaly score function. Interestingly, best performances were achieved when $\beta < 1$ and the performances generally degrade as β increases, resonating with [36] that overly large values of β actually causes a mismatch between $q_\phi(z)$ and $p(z)$ (resulting in inappropriate level of overlap in the latent space). This phenomenon can be further seen in the t-SNE maps of latent embeddings in Fig. 4. When β becomes extremely larger than the appropriate value, the anomalous class becomes en-

tangled with its normal neighboring classes and the boundaries between normal classes become unclear. Theoretically, adding the divergence of the aggregate posterior $q_\phi(\mathbf{z})$ and the desired structured prior $p(\mathbf{z})$ is an effective way to limit the level of overlap when β is too large. However, it is practically challenging to design an appropriate structured prior. Therefore, in our investigation, we focused on exploring the full range of β 's value in β -VAE for the impact of disentanglement to anomaly detection. Since our framework is quite general, it can be easily extended to other unsupervised disentangled representation learning models for anomaly detection.

C. Anomaly Scores with Normalized k -NN Distance in t -SNE Maps

Since the normalised reconstruction error in input space and the k -NN distance in t -SNE maps may have very different magnitudes (as mentioned in Section III-E), some values of α in Equation (2) are close to (but not exactly equal to) 1. The reader may have the intuition that t -SNE is not useful in AnoDM. A simple way, to find out whether large α value is due to magnitude difference or useless of t -SNE, is to replace k -NN distance of a test sample in t -SNE map (Equation (4)) with a k -NN distance normalized by average distance among training samples in the t -SNE map as defined below:

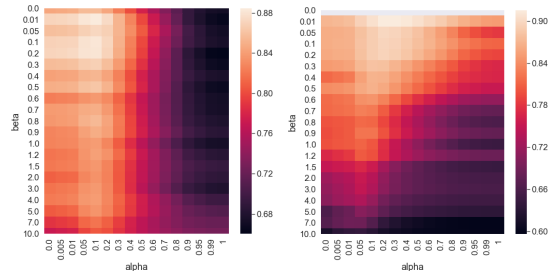
$$\mathcal{N}_{\text{tSNE}}^k(\mathbf{x}_{\text{te}}^{(i)}) \triangleq \frac{\mathcal{D}_{\text{tSNE}}^k(\mathbf{x}_{\text{te}}^{(i)})}{c * \mathcal{D}_{\text{tSNE}}(\mathbf{x}_{\text{tr}})}, \quad (5)$$

$$\mathcal{D}_{\text{tSNE}}(\mathbf{x}_{\text{tr}}) \triangleq \frac{1}{n} \sum_{i,j \in \{1,2,\dots,n\}} \|\mathbf{x}_{\text{tr}}^{(i)} - \mathbf{x}_{\text{tr}}^{(j)}\|_2, \quad (6)$$

where $c > 0$ is normalization hyperparameter (we set it to 0.5 in our experiment); \mathbf{x}_{tr} is a training sample; n is the total number of training samples in a t -SNE map. The heatmaps in Fig. 5 depict performances of AnoDM on Fashion-MNIST (“Dress” class is considered as anomaly) and Arrhythmia (“F” class is treated as anomaly) using normalized k -NN distance in t -SNE maps in combination reconstruction error. By comparing Fig. 5 with Fig. 3, one can see that the optimal values of α shift upper right corner area to the upper left corner area. It thus implies that, the optimal value of α is affected by the magnitude difference, and t -SNE indeed plays an essential role in AnoDM.

D. Impact of Beta to t -SNE Representations

The t -SNE plots in Fig. 4 reflect the impact of β 's value on latent representations in the case of identifying anomalous digit 7 from MNIST. In this example, the best performance (auROC = 0.975) was achieved when $\beta = 0.01$. Clearly, as β increases, all latent clusters become less dense, and more anomalous latent data points move to neighboring clusters. Furthermore, Fig. 4 also corroborates that, even though in t -SNE maps distances between clusters might not reflect global geometry and cluster sizes might not mirror the true sizes [37], using averaged distance from a test sample to its k nearest normal data points represented in t -SNE space to qualify outlieriness still is a very effective way for distinguishing anomalous samples when β is tuned properly.



(a) Fashion-MNIST, “Dress” class. (b) Arrhythmia, “F” class.

Fig. 5: Performances (measured in terms of auROC) of AnoDM evaluated on Fashion-MNIST and Arrhythmia when using normalized k -NN distance in t -SNE maps in combination with reconstruction error. On each dataset, the anomalous class is indicated in the corresponding subcaption while treating the rest classes as normal classes.

E. Evaluation of Anomaly Score Function

In order to better evaluate our anomaly score function, as formulated in Equation (2), we conducted a comprehensive comparison with methods only based on either distance in t -SNE map ($\mathcal{D}_{\text{tSNE}}^k$) or reconstruction-error in raw feature space (\mathcal{D}_{RE}). To see the contribution of t -SNE, it is also compared with the method that calculates nearest neighbor distance directly in latent space of β -VAE. Fig. 6 displays the ROC curves of these four approaches when assuming anomalous classes are respectively 1 (“Trouser/pants”), 3 (“Dress”), 5 (“Sandal”), and 7 (“Sneaker”) on Fashion-MNIST. It is obvious that AnoDM achieves best results among them by taking advantages of both β -VAE reconstruction and t -SNE embedding. The β -VAE reconstruction reflects whether useful information is captured by the model through recovering the input \mathbf{x} ; the t -SNE embedding indicates the disentanglement of latent representations \mathbf{z} . Both measures effectively complement each other. Besides, comparing the auROCs between t -SNE-based and latent-distance-based score functions, one can clearly see that the former dramatically outperforms the latter. Same conclusion can be drawn for MNIST, CIFAR-10, and Arrhythmia. To further show that the optimal values of α are close to 1 (see Fig. 3) is due to magnitude difference rather than less usefulness of t -SNE, we replaced the distance score (Equation (4)) with normalized distance score (Equation (5)) in the weighted final anomaly score function (Equation (2)). We found that the optimal values of α shift to the lower end of the spectrum (see Fig. 5). It implies that t -SNE does play a critical role in our framework.

F. AnoDM for Time-Series

As mentioned in Section III, our method uses a TCN encoder in β -VAE for time-series anomaly detection. Fig. 7 displays the comparison among TCN, CNN and LSTM encoders in the AnoDM framework on Arrhythmia. As a special case of LSTM- β -VAE, LSTM-VAE was presented in [27] for state-of-the-art sequence modelling. For the five classes in Arrhythmia, iteratively one class was treated as anomalous class, while the other classes were used as normal classes. The TCN-encoder-based method outperforms the other two

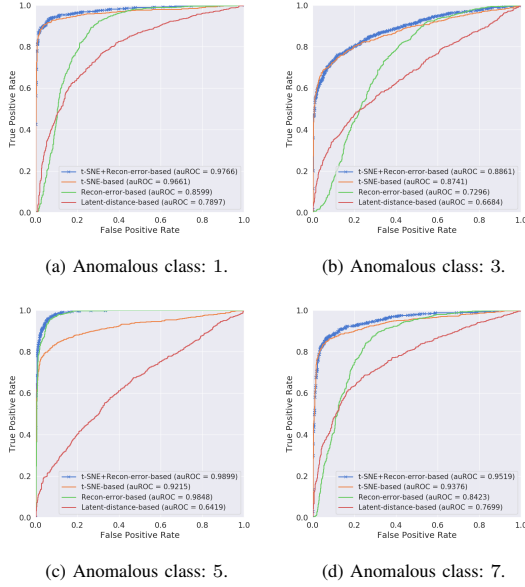


Fig. 6: ROC curves of four methods on Fashion-MNIST. These four examples illustrate the results of using four different anomaly score functions: t-SNE+Recon-error-based ($S_{\beta\text{VAE}+\text{tSNE}} = \alpha\mathcal{D}_{\text{RE}} + (1-\alpha)\mathcal{D}_{\text{tSNE}}^k$), t-SNE-based ($\mathcal{D}_{\text{tSNE}}^k$), Reconstruction-error-based (\mathcal{D}_{RE}) and latent-distance-based (calculating distances in latent space of β -VAE). The α values for these four plots are 0.8, 0.8, 0.95, and 0.8, respectively.

methods significantly in all five cases. Even though the CNN encoder achieved impressive results when detecting anomalous class “S”, “V”, “F” and “Q” respectively, it did not work quite well when class “N” was treated as anomaly. One possible reason might be that comparing with TCN and LSTM, the performance of CNN is more sensitive on the training sample size. Taking the above case as an example, as class 0 (“N”) accounts for over 80% of training data, when considering it as anomaly, normal training data hence become less sufficient for learning β -VAE. Nevertheless, in TCN-based β -VAE, each hidden unit of the last deterministic hidden layer before latent

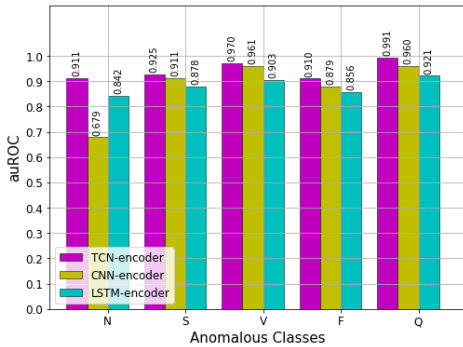
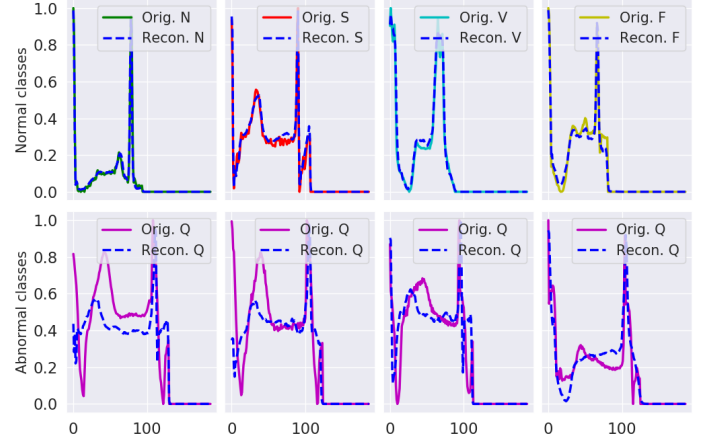


Fig. 7: AuROCs on Arrhythmia.

encoding at the bottleneck is calculated based on much longer sequence dependency, such that it is less sensitive to the limitation of small sample size. Conclusively, as mentioned in [17], TCN should be regarded as a natural starting point for sequence modeling tasks.

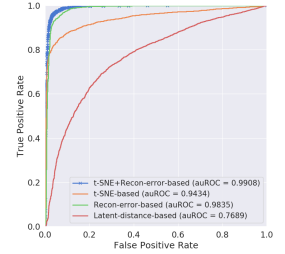
As a case study, Fig. 8a shows the original EGG signals and reconstructed signals by TCN-based β -VAE when considering class “Q” as anomaly. Normal samples can be reconstructed very well, whereas abnormal samples suffer from larger reconstruction errors. Meanwhile, the corresponding t-SNE plot in Fig. 8b displays two distinctive clusters of abnormal samples. The combination of these two measures thus leads to the best performance as seen in Fig. 8c.



(a) Reconstructed signals by β -VAE.



(b) t-SNE plot.



(c) ROC plot.

Fig. 8: Reconstructed signals, t-SNE map, and ROC curves on Arrhythmia with class “Q” as anomaly.

V. CONCLUSIONS

We propose a new methodology which successfully integrates t-SNE with disentangled representation learning for anomaly detection. This approach achieved state-of-the-art performances on both image data (MNIST, Fashion-MNIST and CIFAR-10) and Arrhythmia time-series data. Specifically, best performance is accomplished when $0 < \beta < 1$ for almost all cases involving β -VAE. We also defined an anomaly score function by effectively taking the advantages of both low-dimensional t-SNE embedding and β -VAE reconstruction. Our algorithm demonstrated that t-SNE plays an essential role for measuring abnormality. This research initiates the research on anomaly detection using unsupervised disentangled representation learning and lower-dimensional manifold learning. Besides, our model uses TCN network as encoding architecture for detecting anomalous time-series data and the experimental results convince us that TCN consistently outperforms CNN

and LSTM. As a proof of concept, our current framework automatically inherits advantages of deep learning to address anomaly detection’s issues in representability and scalability as discussed in the beginning of this paper. The extension of our framework to multimodal data is straightforward. It is also possible that a neural t-SNE component could be designed and integrated into the learning of β -VAE to achieve real-time efficiency. Other new well-performing manifold learning methods, such as UMAP [44] which is faster and keeps global topologies, could be employed as replacement of t-SNE.

APPENDIX

A. Beta-VAE and Beyond for Disentangled Representation Learning

[13] proposed a novel deep generative model, named β -VAE, a modification of VAE by introducing an adjustable hyperparameter β to learn an interpretable disentangled representation of the data generative latent factors. Specifically, β functions as a controller to trade off between the extent of learning constraints and reconstruction accuracy. The constraints impose a limit on the capacity of the latent information channel and an emphasis on learning statistically independent latent factors. [13] demonstrated that β -VAE with appropriately tuned β ($\beta > 1$) qualitatively outperforms VAE ($\beta = 1$) as well as state of the art unsupervised (InfoGAN) and semi-supervised (DC-IGN) approaches to disentangled factor learning on a variety of datasets (celebA, faces and chairs).

[13] assumed that an image \mathbf{x} is generated by the true world simulator using ground truth data generative factors: $p(\mathbf{x}|\mathbf{v}, \mathbf{w}) = \text{Sim}(\mathbf{v}, \mathbf{w})$, where \mathbf{v} is set of conditionally independent factors and \mathbf{w} is set of conditionally dependent factors. Therefore, the joint distribution of the data \mathbf{x} and a set of generative latent factors \mathbf{z} is: $p(\mathbf{x}|\mathbf{z}) \approx p(\mathbf{x}|\mathbf{v}, \mathbf{w}) = \text{Sim}(\mathbf{v}, \mathbf{w})$. The aim of this generative model is then to ensure that the inferred latent factors from $q_\phi(\mathbf{z}|\mathbf{x})$ capture the generative factors \mathbf{v} in a disentangled manner. The conditionally dependent data generative factors \mathbf{w} can remain entangled in a separate subset of \mathbf{z} that is not entangled with \mathbf{v} . Considering the prior $p(\mathbf{z})$ is set to be an isotropic unit Gaussian $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$, a constraint δ is introduced to encourage the matching between $q_\phi(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{z})$ such that the disentangling property in the inferred $q_\phi(\mathbf{z}|\mathbf{x})$ can be realized. Following the same incentive as in VAE: maximizing the probability of generating real data, while minimizing the distance between the generative and approximate posterior distributions, as formulated below

$$\max_{\phi, \theta} \mathbb{E}_{\mathbf{x} \sim \mathbb{X}} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]] \quad (7)$$

$$\text{s.t. } D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) < \delta, \quad (8)$$

where $\mathbb{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ is the training data set. The objective function to be maximized in β -VAE is thus defined as:

$$\mathcal{L}_\beta(\theta, \phi) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})), \quad (9)$$

where the Lagrangian multiplier β is the regularisation coefficient that constrains the capacity of the latent information channel \mathbf{z} and puts implicit independence pressure on the learnt posterior due to the isotropic nature of the Gaussian prior $p(\mathbf{z})$. When $\beta = 1$, β -VAE corresponds to the original VAE formulation of [12]. When $\beta > 1$, it applies stronger constraint which limits the capacity of \mathbf{z} and encourages the model to learn the most efficient representation of the data. Theoretically, a higher β encourages more efficient latent encoding and further encourages the disentanglement. However, a higher β may lead to poorer reconstructions due to the loss of high frequency details when passing through a constrained latent bottleneck.

Mathieu et al. [36] argued that overly large β is not universally beneficial for disentanglement, Since this in turn causes a mismatch between marginal posterior $q_\phi(\mathbf{z})$ and the prior $p(\mathbf{z})$. Thus they proposed a generalization of disentanglement in VAE by explicitly separating such a decomposition as two tasks: a) the latent encoding of data should achieved an appropriate level of non-negligible overlap in aggregate encoding $q_\phi(\mathbf{z})$, and b) the aggregate encoding of data $q_\phi(\mathbf{z})$ should match the prior $p(\mathbf{z})$ which demonstrates the desired dependency structure between latent variables. [36] developed a new objective that incorporates both a) and b) by introducing an additional divergence term $\mathbb{D}(q_\phi(\mathbf{z}), p(\mathbf{z}))$.

$$\begin{aligned} \mathcal{L}_{\alpha, \beta}(\mathbf{x}) = & \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \\ & - \alpha \mathbb{D}(q_\phi(\mathbf{z}), p(\mathbf{z})). \end{aligned} \quad (10)$$

By appropriately setting β and α , it allows direct control over the level of overlap and the regularization between the marginal posterior and the prior. However, a practical challenge of this method is how to define a proper structured prior when the structure of real hidden factors is poorly known. For this reason, our computational experiment in this paper is based on [14]’s β -VAE.

REFERENCES

- [1] D. M. Tax and R. P. Duin, “Support vector domain description,” *Pattern Recognition Letters*, vol. 20, pp. 1191–1199, 1999.
- [2] N. Gornitz, M. Kloft, K. Rieck, and U. Brefeld, “Toward supervised anomaly detection,” *Journal of Artificial Intelligence Research*, vol. 43, pp. 235–262, 2013.
- [3] Y. LeCun, Y. Bengio, and G. E. Hinton, “Deep learning,” *Nature*, vol. 521, pp. 436–444, 2015.
- [4] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley, “Stochastic variational inference,” *Journal of Machine Learning Research*, vol. 14, pp. 1303–1347, 2013.
- [5] C. Zhang, J. Butepage, H. Kjellstrom, and S. Mandt, “Advances in variational inference,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 8, pp. 2008 – 2026, 2018.
- [6] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in PyTorch,” in *Neural Information Processing Systems Autodiff Workshop*, 2017.
- [7] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng, “TensorFlow: A system for large-scale machine learning,” in *12th USENIX Symposium on Operating Systems Design and*

TABLE III: Architecture of the β -VAE used in AnODM.

Dataset	Optimizer	Architecture	
MNIST	Adam ($1e^{-3}$)	Input	784 (flattened $28 \times 28 \times 1$).
		Encoder	Conv2d (filters 32, kernel_size 3, strides 2), Conv2d (filters 32, kernel_size 3, strides 2), FC (128), FC (128), ReLU activation.
		Decoder	FC (128), FC (1568), Deconv2d (filters 32, kernel_size 3, strides 2), Deconv2d (filters 1, kernel_size 3, strides 2), ReLU activation.
Fashion-MNIST	Adam ($1e^{-3}$)	Input	784 (flattened $28 \times 28 \times 1$).
		Encoder	Conv2d (filters 32, kernel_size 3, strides 2), Conv2d (filters 32, kernel_size 3, strides 2), FC (128), FC (128), ReLU activation.
		Decoder	FC (128), FC (1568), Deconv2d (filters 32, kernel_size 3, strides 2), Deconv2d (filters 1, kernel_size 3, strides 2), ReLU activation.
CIFAR-10	Adam ($1e^{-3}$)	Input	3072 (flattened $32 \times 32 \times 3$).
		Encoder	Conv2d (filters 64, kernel_size 3, strides 2), BN, Conv2d (filters 128, kernel_size 5, strides 2), BN, Conv2d (filters 128, kernel_size 5, strides 2), BN, FC (256), Dropout (0.9), BN, FC (256), Dropout (0.9), BN, ReLU activation.
		Decoder	FC (256), BN, FC (256), BN, FC (2048), BN, Deconv2d (filters 128, kernel_size 5, strides 2), BN, Deconv2d (filters 64, kernel_size 5, strides 2), BN, Deconv2d (filters 3, kernel_size 3, strides 2), ReLU activation.
Small-Norb	Adam ($1e^{-3}$)	Input	1024 (flattened $32 \times 32 \times 1$).
		Encoder	Conv2d (filters 32, kernel_size 3, strides 2), BN, Conv2d (filters 64, kernel_size 3, strides 2), BN, Conv2d (filters 128, kernel_size 3, strides 2), BN, FC (256), Dropout (0.7), BN, FC (256), Dropout (0.7), BN, ReLU activation.
		Decoder	FC (256), BN, FC (256), BN, FC (2048), BN, Deconv2d (filters 128, kernel_size 5, strides 2), Deconv2d (filters 64, kernel_size 5, strides 2), Deconv2d (filters 3, kernel_size 3, strides 2), ReLU activation.
Arrhythmia	Adam ($1e^{-3}$)	Input	187.
		Encoder	TCN (filters 64, kernel_size 4, stacks 1, dilations [1, 2, 4, 8, 16, 32], padding "causal", use_skip_connections "True", dropout_rate 0.05)
		Decoder	FC (256), FC (384), FC (2048), BN, Deconv1d (filters 64, kernel_size 3, strides 2), Deconv2d (filters 64, kernel_size 3, strides 2), Deconv2d (filters 64, kernel_size 3, strides 2), Deconv2d (filters 32, kernel_size 3, strides 2), Deconv2d (filters 32, kernel_size 3, strides 2), Deconv2d (filters 1, kernel_size 3, strides 2), ReLU activation.

Implementation (OSDI 16), 2016, pp. 265–283. [Online]. Available: <https://www.usenix.org/system/files/conference/osdi16/osdi16-abadi.pdf>

- [8] Y. Li, F.-X. Wu, and A. Ngom, "A review on machine learning principles for multi-view biological data integration," *Briefings in Bioinformatics*, vol. 19, no. 2, pp. 325–340, 2018.
- [9] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [10] E. Nalisnick, A. Matsukawa, Y. W. Teh, D. Gorur, and B. Lakshminarayanan, "Do deep generative models know what they don't know?" in *International Conference on Learning Representations*, 2019.
- [11] J. An and S. Cho, "Variational autoencoder based anomaly detection using reconstruction probability," Data Mining Center, Seoul National University, Seoul, South Korea, Tech. Rep., 2015.
- [12] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," in *International Conference on Learning Representations*, 2014.
- [13] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "beta-VAE: Learning basic visual concepts with a constrained variational framework," *International Conference on Learning Representations*, vol. 2, no. 5, p. 6, April 2017.
- [14] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, and A. Lerchner, "Understanding disentangling in β -VAE," *ArXiv*, p. arXiv:1804.03599v1, April 2018.
- [15] M. D. Hoffman, C. Riquelme, and M. J. Johnson, "The β -VAE's implicit priors," in *Neural Information Processing Systems 2017 Workshop Bayesian Deep Learning*, vol. 2, 2017.
- [16] L. van der Maaten and G. E. Hinton, "Visualizing data using t-SNE," *Journal of Machine Learning Research*, vol. 9, pp. 2579 – 2605, 2008.
- [17] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *ArXiv*, p. arXiv:1803.01271v2, 2018.
- [18] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, pp. 1735–80, December 1997.
- [19] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural Computation*, vol. 18, pp. 1527–1554, 2006.
- [20] Y. Li and X. Zhu, "Exploring Helmholtz machine and deep belief net in the exponential family perspective," in *International Conference on Machine Learning 2018 Workshop on Theoretical Foundations and Applications of Deep Generative Models*, July 2018.
- [21] —, "Exponential family restricted Boltzmann machines and annealed importance sampling," in *International Joint Conference on Neural Networks*, July 2018, pp. 39–48.
- [22] J. T. A. Andrews, E. J. Morton, and L. D. Griffin, "Detecting anomalous data using auto-encoders," *International Journal of Machine Learning and Computing*, vol. 6, no. 1, pp. 21–26, 2016.
- [23] S. M. Erfani, S. Rajasegarar, S. Karunasekera, and C. Leckie, "High-dimensional and large-scale anomaly detection using a linear one-class SVM with deep learning," *Pattern Recognition*, vol. 58, pp. 121–134, 2016.
- [24] T. Schlegl, P. Seebock, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs, "Unsupervised anomaly detection with generative adversarial networks to guide marker discovery," in *International Conference on Information Processing in Medical Imaging*. Springer, 2017, pp. 146–157.
- [25] L. Deecke, R. Vandermeulen, L. Ruff, S. Mandt, and M. Kloft, "Image anomaly detection with generative adversarial networks," in *Machine Learning and Knowledge Discovery in Databases*. Cham: Springer International Publishing, 2019, pp. 3–17.
- [26] D. Li, D. Chen, J. Goh, and S. Ng, "Anomaly detection with generative adversarial networks for multivariate time series," *ArXiv*, p. arXiv:1809.04758, 2018.
- [27] D. Park, Y. Hoshi, and C. C. Kemp, "A multimodal anomaly detector for robot-assisted feeding using an LSTM-based variational autoencoder," *IEEE Robotics and Automation Letters*, vol. 3, no. 3, 2017.
- [28] K. Cho, B. van Merriënboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," *ArXiv*, p. arXiv:1409.1259, 2014.
- [29] R. Jozefowicz, W. Zaremba, and I. Sutskever, "An empirical exploration of recurrent network architectures," in *International Conference on Machine Learning*, 2015, pp. 2342–2350.
- [30] K. Greff, R. K. Srivastava, J. Koutník, B. R. Steunebrink, and J. Schmidhuber, "LSTM: A search space odyssey," *ArXiv*, p. arXiv:1503.04069, 2015.
- [31] G. Melis, C. Dyer, and P. Blunsom, "On the state of the art of evaluation in neural language models," *ArXiv*, p. arXiv:1707.05589, 2017.
- [32] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, pp. 541–551, Dec. 1989.
- [33] J. Gehring, M. Auli, D. Grangier, and Y. N. Dauphin, "A convolutional encoder model for neural machine translation," *ArXiv*, p. arXiv:1611.02344, 2016.
- [34] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin, "Convolutional sequence to sequence learning," in *International Conference on Machine Learning*, vol. 70, 2017, pp. 1243–1252.
- [35] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier, "Language modeling with gated convolutional networks," *ArXiv*, p. arXiv:1612.08083, 2016.
- [36] E. Mathieu, T. Rainforth, N. Siddharth, and Y. W. Teh, "Disentangling disentanglement in variational autoencoders," in *ICML*, vol. 97. Long Beach, California, USA: PMLR, Jun 2019, pp. 4402–4412.
- [37] M. Wattenberg, F. Viégas, and I. Johnson, "How to use t-SNE effectively," *Distill*, 2016. [Online]. Available: <http://distill.pub/2016/misread-tsne>
- [38] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [39] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms," *ArXiv*, p. arXiv:1708.07747, 2017.
- [40] Y. LeCun, F. J. Huang, and L. Bottou, "Learning methods for generic object recognition with invariance to pose and lighting," in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2004.
- [41] A. Krizhevsky, "Learning multiple layers of features from tiny images," Department of Computer Science, University of Toronto, Toronto, Canada, Tech. Rep., 2009.
- [42] S. Fazeli, "Ecg heartbeat categorization dataset," 2018, accessed: 2019-08-20. [Online]. Available: <https://www.kaggle.com/shayanfazeli/heartbeat>
- [43] X. Li, I. Kiringa, T. Yeap, X. Zhu, and Y. Li, "Exploring deep anomaly detection methods based on capsule net," *International Conference on Machine Learning 2019 Workshop on Uncertainty and Robustness in Deep Learning*, p. arXiv:1907.06312, 2019.
- [44] L. McInnes, J. Healy, and J. Melville, "UMAP: Uniform manifold approximation and projection for dimension reduction," *ArXiv*, p. arXiv:1802.03426, 2018.