

Unsupervised Anomaly Detection in Time Series Data using Deep Learning

Integrated Master in Electrical and Computer Engineering



João Pereira
23rd November, 2018

Introduction & Motivation

Anomaly detection is about finding patterns in data that do not conform to *expected* or *normal* behaviour.

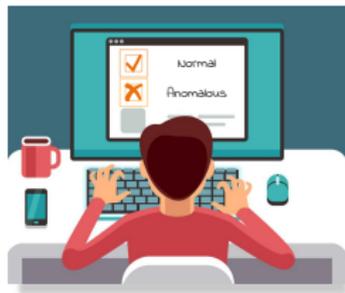


Main Challenges

- ▶ Most data in the world are **unlabelled**

$$\text{Dataset } \mathcal{D} = \left\{ \left(\mathbf{x}^{(i)}, \mathbf{y}^{*(i)} \right) \right\}_{i=1}^N \quad \text{anomaly labels}$$

- ▶ Annotating large datasets is difficult, time-consuming and expensive



- ▶ Time series have temporal structure/dependencies

$$\mathbf{x} = \left(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T \right) , \quad \mathbf{x}_t \in \mathbb{R}^{d_x}$$

We would like:

- ▶ Unsupervised: no need for anomaly labels;
- ▶ Suitable for sequential data (e.g., time series);
- ▶ General;
- ▶ Scalable & efficient, allowing real-time detection.

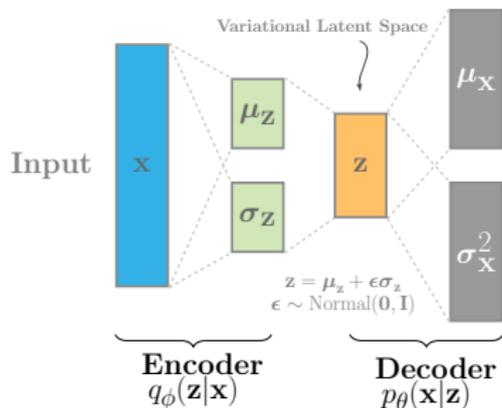
We would like:

- ▶ Unsupervised: no need for anomaly labels;
- ▶ Suitable for sequential data (e.g., time series);
- ▶ General;
- ▶ Scalable & efficient, allowing real-time detection.

How to design such a model?

The Principle in a Nutshell

- ▶ Train a **Variational Autoencoder**¹² to reconstruct input data with mostly normal patterns;



- ▶ At test time, it reconstructs well *normal* data, while it fails to reconstruct *anomalous* data;
- ▶ The quality of the reconstructions and the representations are used to compute anomaly scores.

¹Kingma & Welling, **Auto-Encoding Variational Bayes**, ICLR'14

²Rezende *et al.*, **Stochastic Backpropagation and Approximate Inference in Deep Generative Models**, ICML'14

Proposed Approach

Representation
Learning

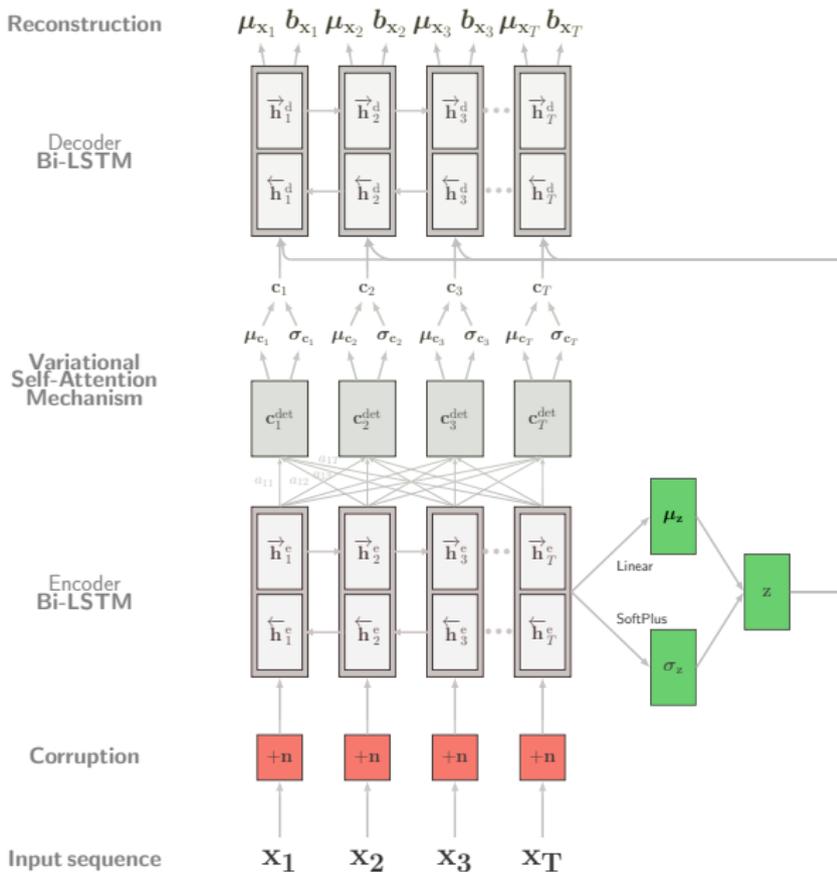
Detection

Proposed Approach

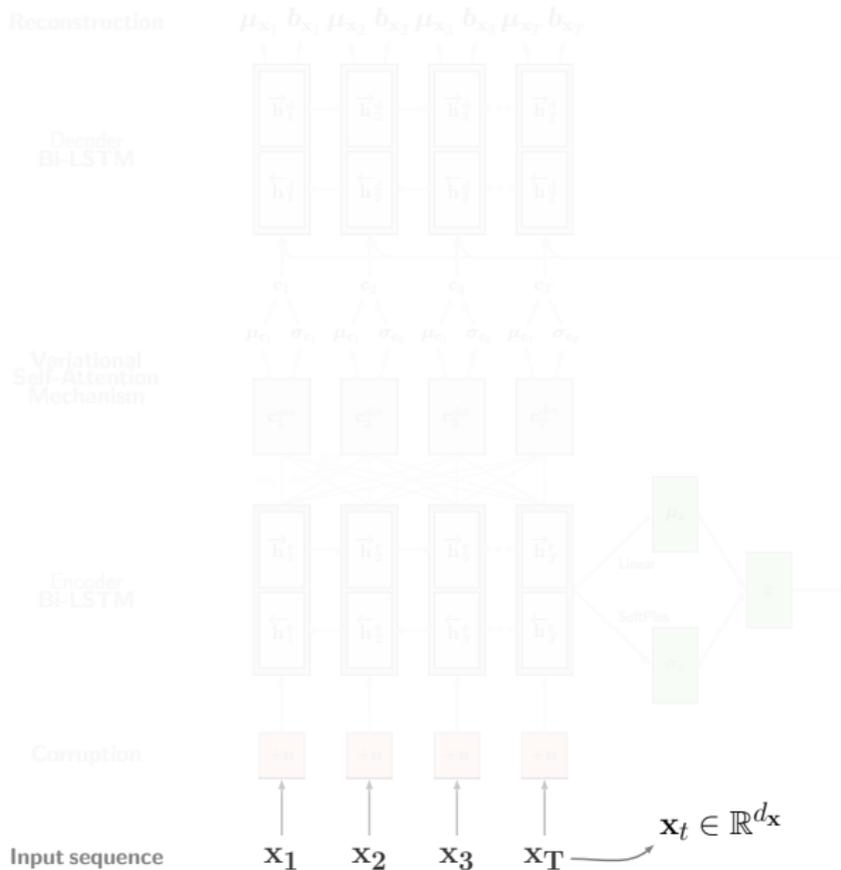
Representation
Learning

Detection

Representation Learning



Representation Learning



Representation Learning

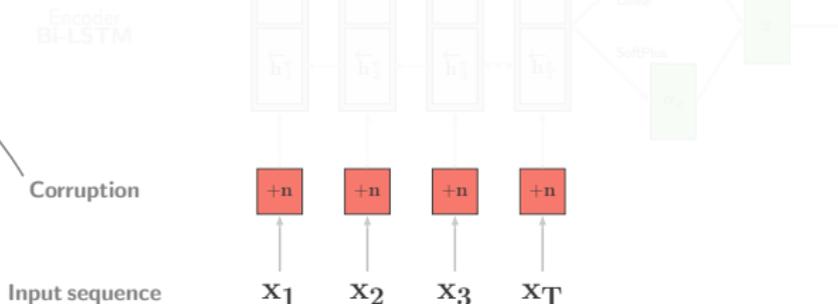
Denosing Autoencoding Criterion

Corruption process: additive Gaussian noise

$$p(\tilde{\mathbf{x}}|\mathbf{x}) = \mathbf{x} + \mathbf{n} \quad , \quad \mathbf{n} \sim \text{Normal}(\mathbf{0}, \sigma_{\mathbf{n}}^2 \mathbf{I})$$

Vincent *et al.*, **Extracting and Composing Robust Features with Denosing Autoencoders**, ICML'08

Bengio *et al.*, **Denosing Criterion for Variational Auto-Encoding Framework**, ICLR'15



Representation Learning

Reconstruction $\mu_{x_1} \ b_{x_1} \ \mu_{x_2} \ b_{x_2} \ \mu_{x_3} \ b_{x_3} \ \mu_{x_T} \ b_{x_T}$

Learning temporal dependencies

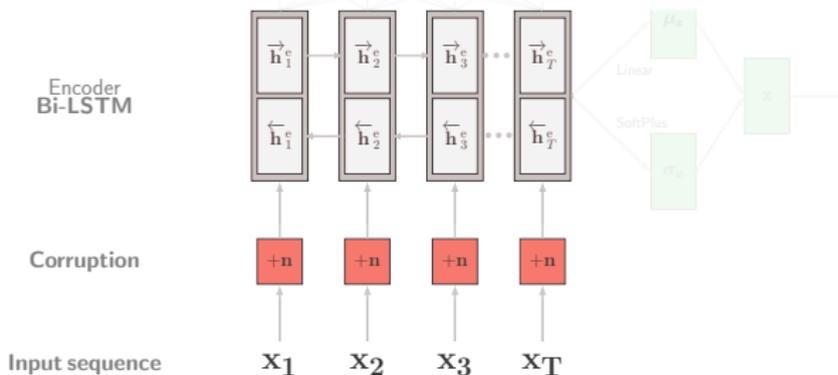
Bidirectional Long-Short Term Memory network

$$\mathbf{h}_t = \left[\vec{\mathbf{h}}_t; \overleftarrow{\mathbf{h}}_t \right]$$

- ▶ 256 units, 128 in each direction
- ▶ Sparse regularization, $\Omega(\mathbf{z}) = \lambda \sum_{i=1}^{d_z} |z_i|$

Hochreiter *et al.*, **Long-Short Term Memory**, Neural Computation'97

Graves *et al.*, **Bidirectional LSTM Networks for Improved Phoneme Classification and Recognition**, ICANN'05



Representation Learning

Variational Latent Space

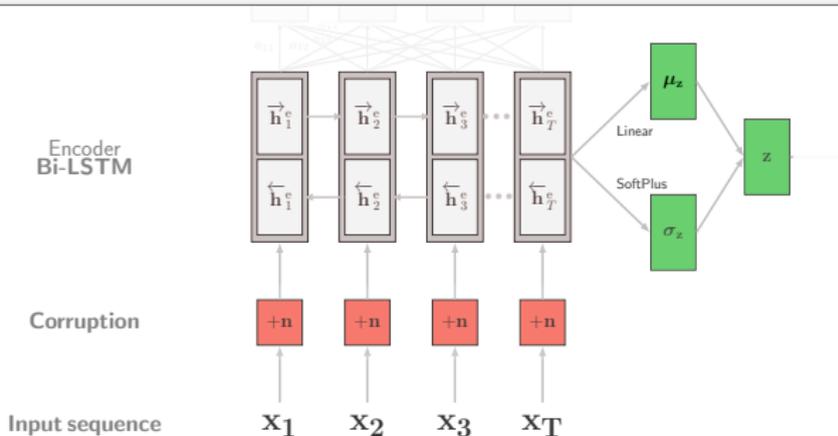
Variational parameters derived using neural networks

$$(\mu_z, \sigma_z) = \text{Encoder}(\mathbf{x})$$

Sample from the approximate posterior $q_\phi(\mathbf{z}|\mathbf{x})$

$$\mathbf{z} = \mu_z + \sigma_z \odot \epsilon \quad \epsilon \sim \text{Normal}(\mathbf{0}, \mathbf{I})$$

Kingma & Welling, **Auto-Encoding Variational Bayes**, ICLR, 2014



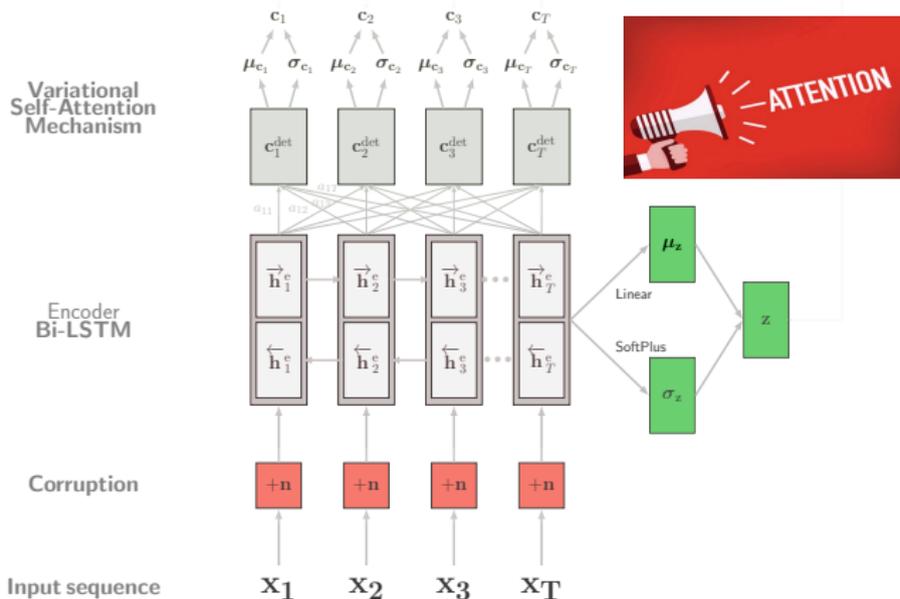
Representation Learning

Introducing Variational Self-Attention

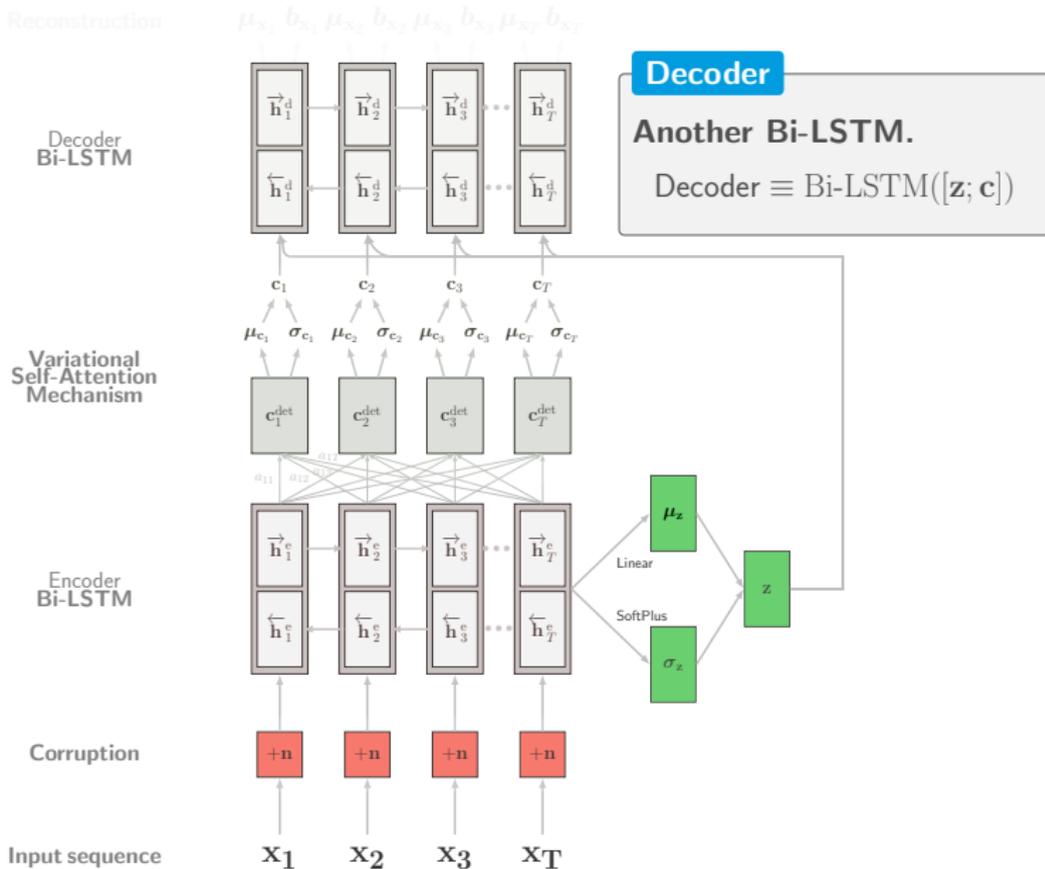
A combination of self-attention and variational inference.

$$\mathbf{c}_t^{\text{det}} = \sum_{j=1}^T a_{tj} \mathbf{h}_j \quad (\boldsymbol{\mu}_{\mathbf{c}_t}, \boldsymbol{\sigma}_{\mathbf{c}_t}) = \text{NN}(\mathbf{c}_t^{\text{det}}), \quad \mathbf{c}_t \sim \text{Normal}(\boldsymbol{\mu}_{\mathbf{c}_t}, \boldsymbol{\sigma}_{\mathbf{c}_t}^2 \mathbf{I})$$

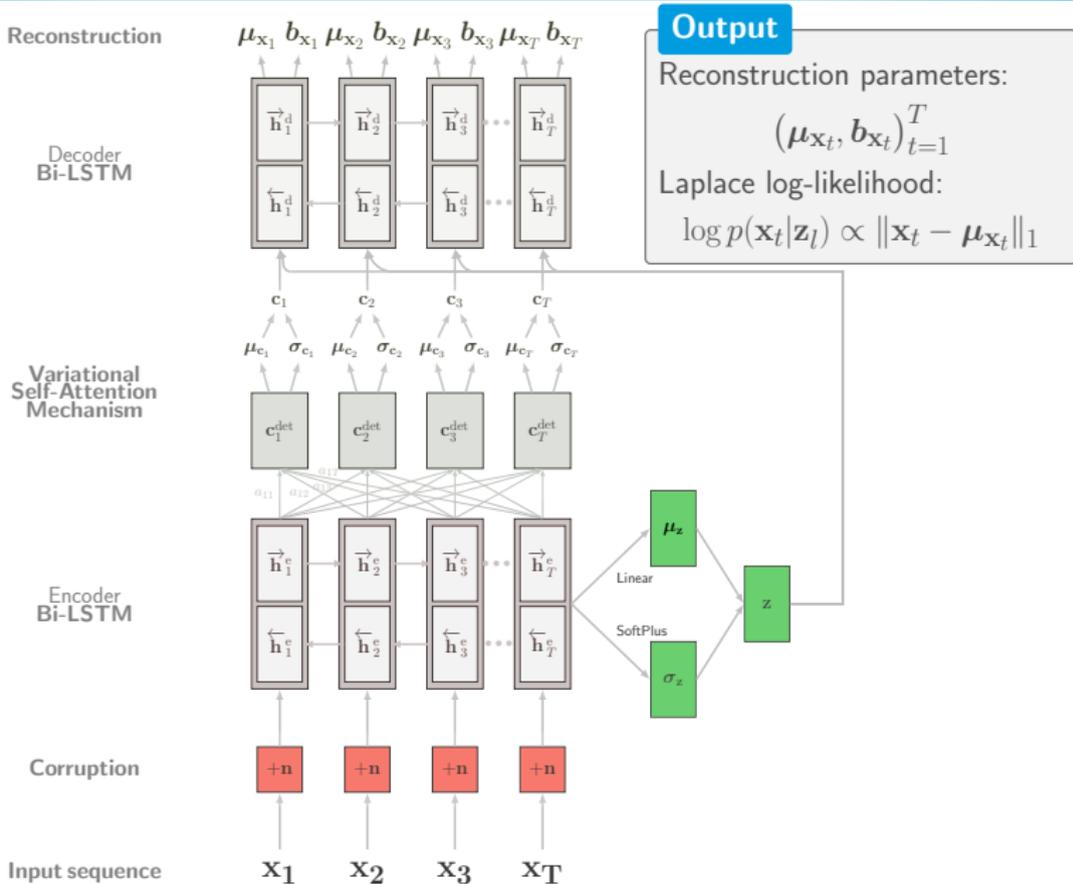
Vaswani *et al.*, Attention is All You Need, NIPS'17



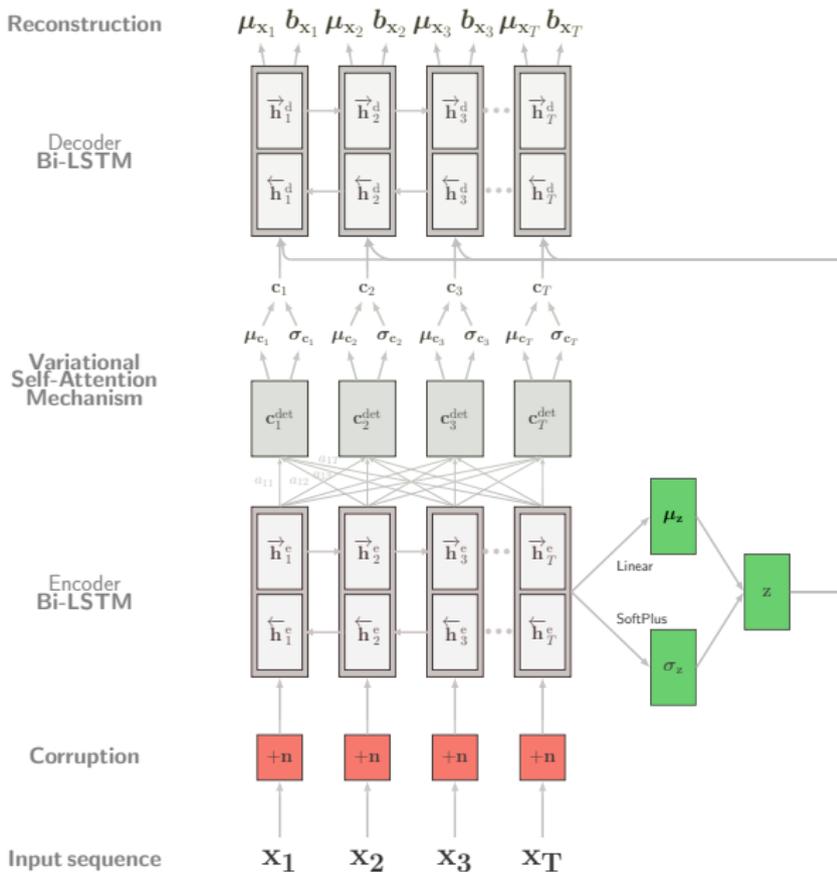
Representation Learning



Representation Learning



Representation Learning



Loss Function

$$\begin{aligned}\mathcal{L}(\theta, \phi; \mathbf{x}^{(n)}) &= -\mathbb{E}_{\mathbf{z} \sim \tilde{q}_\phi(\mathbf{z}|\mathbf{x}^{(n)}), \mathbf{c}_t \sim \tilde{q}_\phi^a(\mathbf{c}_t|\mathbf{x}^{(n)})} \left[\log p_\theta(\mathbf{x}^{(n)}|\mathbf{z}, \mathbf{c}) \right] \\ &+ \lambda_{\text{KL}} \left[\mathcal{D}_{\text{KL}}\left(\tilde{q}_\phi(\mathbf{z}|\mathbf{x}^{(n)})\|p_\theta(\mathbf{z})\right) + \eta \sum_{t=1}^T \mathcal{D}_{\text{KL}}\left(\tilde{q}_\phi^a(\mathbf{c}_t|\mathbf{x}^{(n)})\|p_\theta(\mathbf{c}_t)\right) \right]\end{aligned}$$

\mathcal{D}_{KL} denotes the Kullback-Leibler Divergence

$$\begin{aligned} \mathcal{L}(\theta, \phi; \mathbf{x}^{(n)}) = & \overbrace{-\mathbb{E}_{\mathbf{z} \sim \tilde{q}_\phi(\mathbf{z}|\mathbf{x}^{(n)}), \mathbf{c}_t \sim \tilde{q}_\phi^a(\mathbf{c}_t|\mathbf{x}^{(n)})} \left[\log p_\theta(\mathbf{x}^{(n)}|\mathbf{z}, \mathbf{c}) \right]}^{\text{Reconstruction Term}} \\ & + \lambda_{\text{KL}} \left[\underbrace{\mathcal{D}_{\text{KL}}\left(\tilde{q}_\phi(\mathbf{z}|\mathbf{x}^{(n)})\|p_\theta(\mathbf{z})\right)}_{\text{Latent Space KL loss}} + \eta \sum_{t=1}^T \underbrace{\mathcal{D}_{\text{KL}}\left(\tilde{q}_\phi^a(\mathbf{c}_t|\mathbf{x}^{(n)})\|p_\theta(\mathbf{c}_t)\right)}_{\text{Attention KL loss}} \right] \end{aligned}$$

\mathcal{D}_{KL} denotes the Kullback-Leibler Divergence

Optimization & Regularization

- ▶ About 270k parameters to optimize
- ▶ *AMS-Grad* optimizer³
- ▶ Xavier weight initialization⁴
- ▶ Denoising autoencoding criterion⁵
- ▶ Sparse regularization in the encoder Bi-LSTM⁶
- ▶ KL cost annealing⁷
- ▶ Gradient clipping⁸

Training executed on a single GPU (NVIDIA GTX 1080 TI)

³Reddi, Kale & Kumar, *On the Convergence of Adam and Beyond*, ICLR'18

⁴Bengio *et al.*, *Understanding the Difficulty of Training Deep Feedforward Neural Networks*, AISTATS'10

⁵Bengio *et al.*, *Denoising Criterion for Variational Auto-Encoding Framework*, AAAI'17

⁶Arpit *et al.*, *Why Regularized Auto-Encoders Learn Sparse Representation?*, ICML'16

⁷Bowman, Vinyals *et al.*, *Generating Sentences from a Continuous Space*, SIGNLL'16

⁸Bengio *et al.*, *On the Difficulty of Training Recurrent Neural Networks*, ICML'13

Proposed Approach

Representation
Learning

Detection

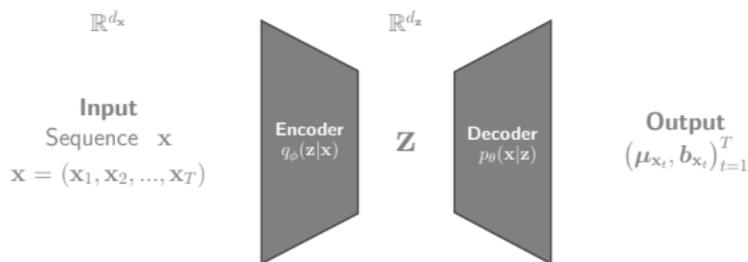
Proposed Approach

Representation
Learning

Detection

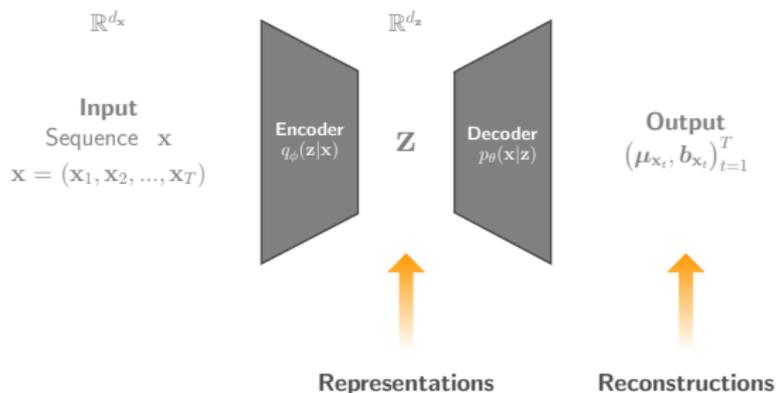
One model, two detection strategies

The model provides two products: **representations** in the z -space and the **reconstruction parameters** in the x -space



One model, two detection strategies

The model provides two products: **representations** in the \mathbf{z} -space and the **reconstruction parameters** in the \mathbf{x} -space

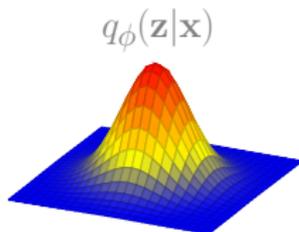


- ▶ Reconstruction-based detection
- ▶ Latent space-based detection

Reconstruction-based Detection

$$\text{Reconstruction Error} = \mathbb{E}_{\mathbf{z}_l \sim q_\phi(\mathbf{z}|\mathbf{x})} \left[\left\| \mathbf{x} - \mathbb{E}[p_\theta(\mathbf{x}|\mathbf{z}_l)] \right\|_1 \right]$$

$$\text{Reconstruction Probability} = \mathbb{E}_{\mathbf{z}_l \sim q_\phi(\mathbf{z}|\mathbf{x})} \left[\log p(\mathbf{x}|\mathbf{z}_l) \right]$$

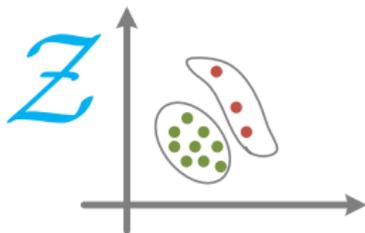


$$\mathbf{z}_l \sim \text{Normal}(\boldsymbol{\mu}_z, \sigma_z^2 \mathbf{I})$$

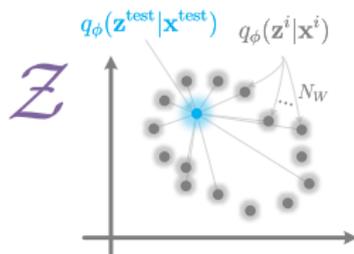
Latent Space Detection

Based on the representations in the z-space.

► Clustering

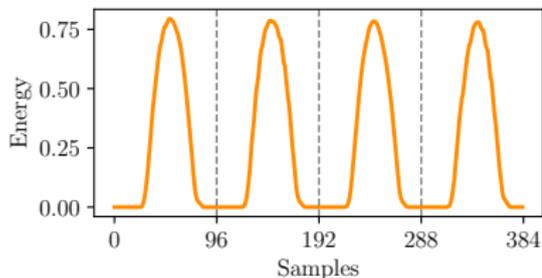


► Wasserstein Metric (W)



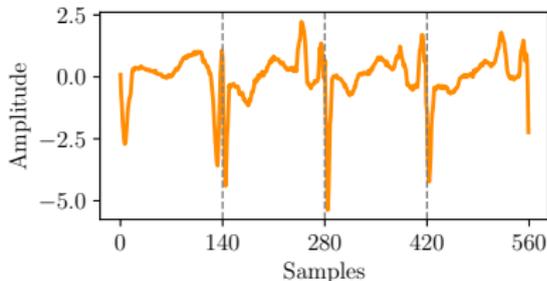
$$\text{score}(\mathbf{z}^{\text{test}}) = \text{median}\{W(\mathbf{z}^{\text{test}}, \mathbf{z}^i)^2\}_{i=1}^{N_W}$$

Solar PV Generation



- ▶ Provided by **c|side**
- ▶ Recorded every 15min (96 samples/day)
- ▶ Daily seasonality
- ▶ Unlabelled

Electrocardiogram



- ▶ Available in the UCR Time Series Classification Archive ECG5000 [Keogh *et al.*, 2015]
- ▶ One heartbeat \approx 140 samples
- ▶ 5000 sequences
- ▶ Labelled, 5 classes annotated

Dataset

Energy



ECG5000



Dataset

Energy



ECG5000



z-space in 2D ($\mathcal{X}_{\text{train}}^{\text{normal}}$)

$T = 12$ (< 96)

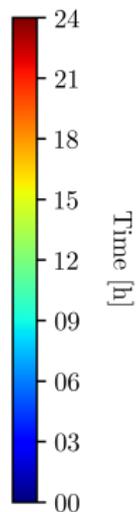
$d_z = 3$

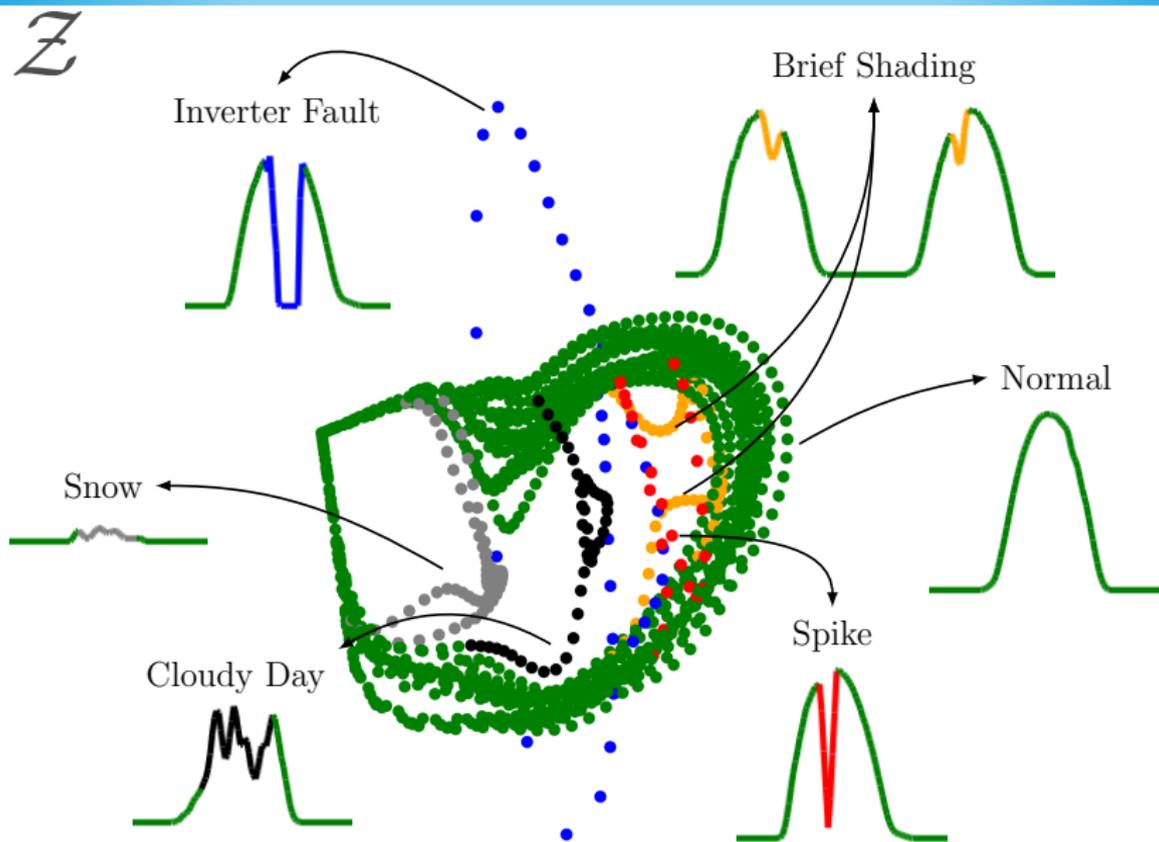
online mode

t-SNE



PCA



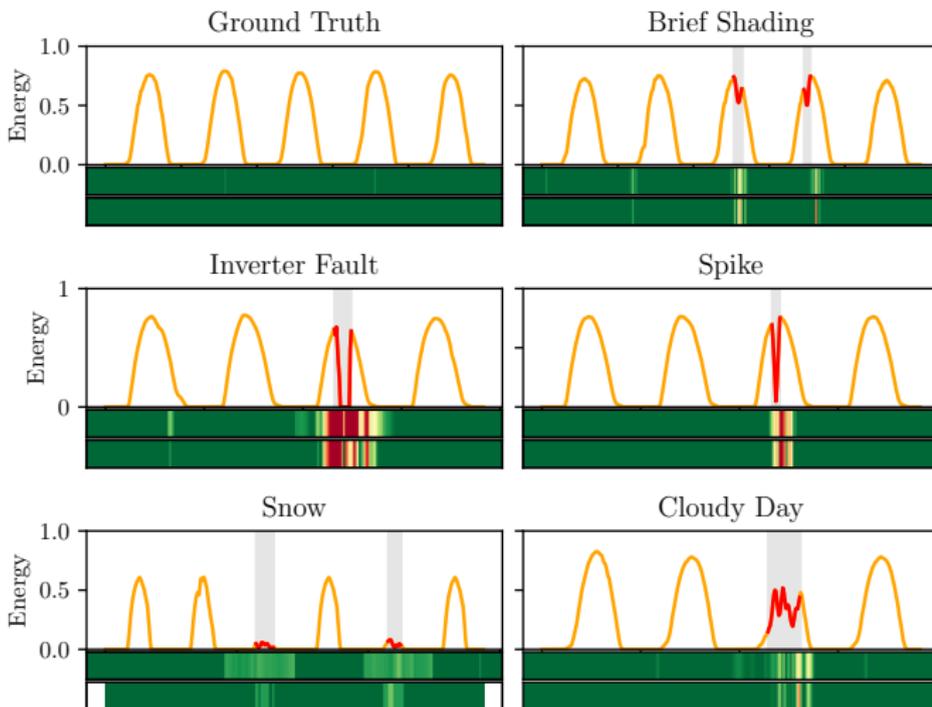


Reconstruction Error

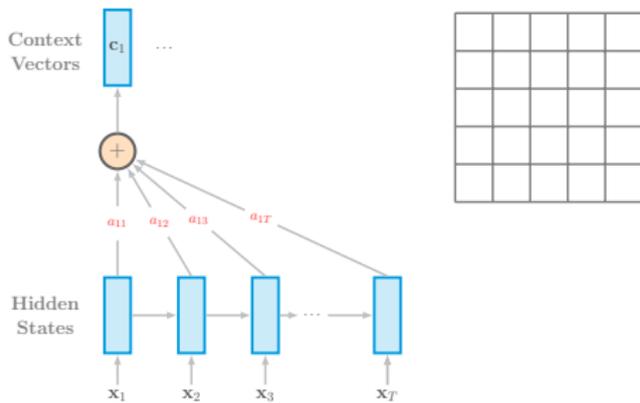
(top bar)

Reconstruction Probability

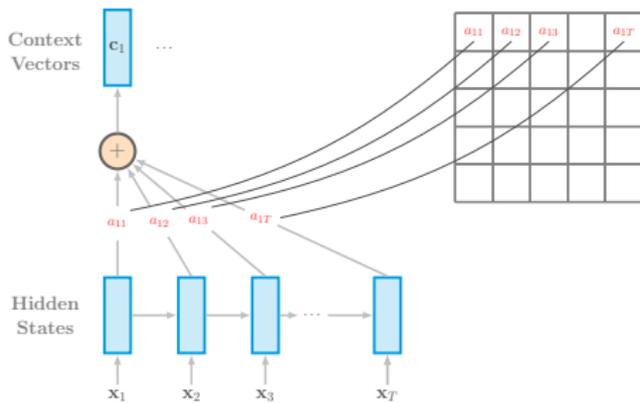
(bottom bar)



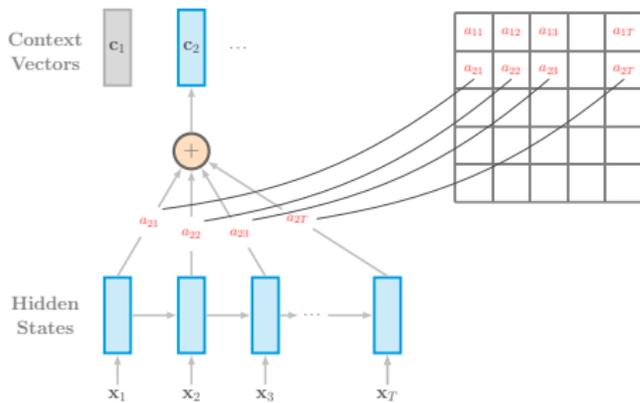
Attention Map



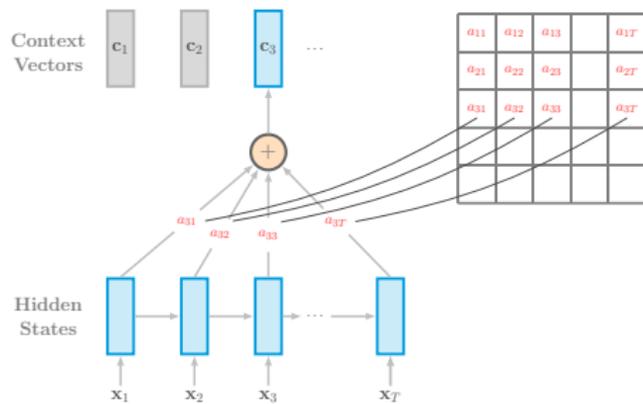
Attention Map



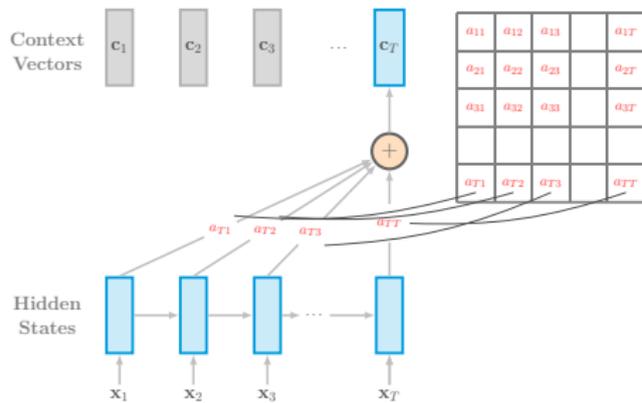
Attention Map



Attention Map

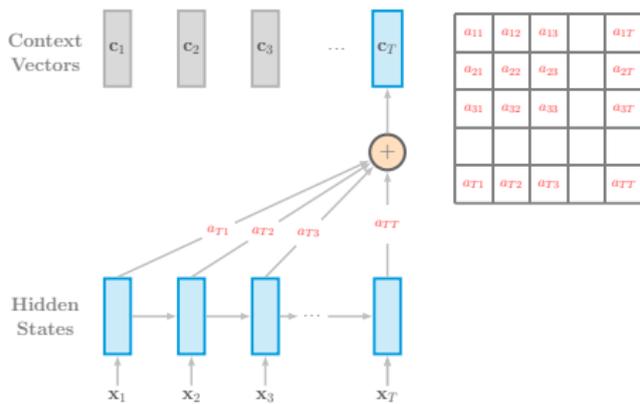


Attention Map

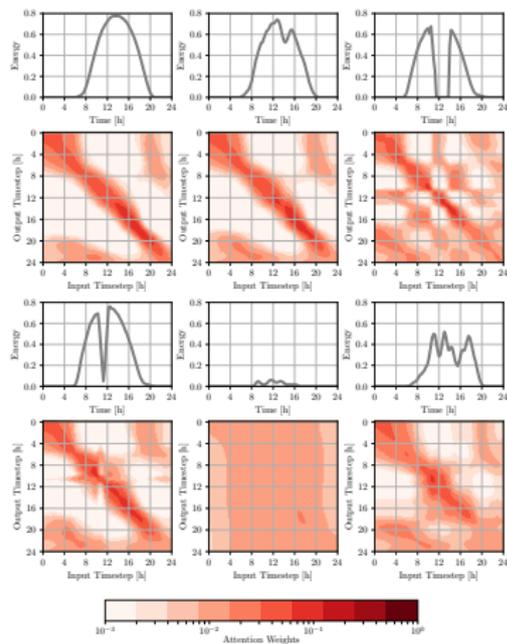


Examples

Attention Map



Examples



Dataset

Energy



ECG5000



Dataset

Energy

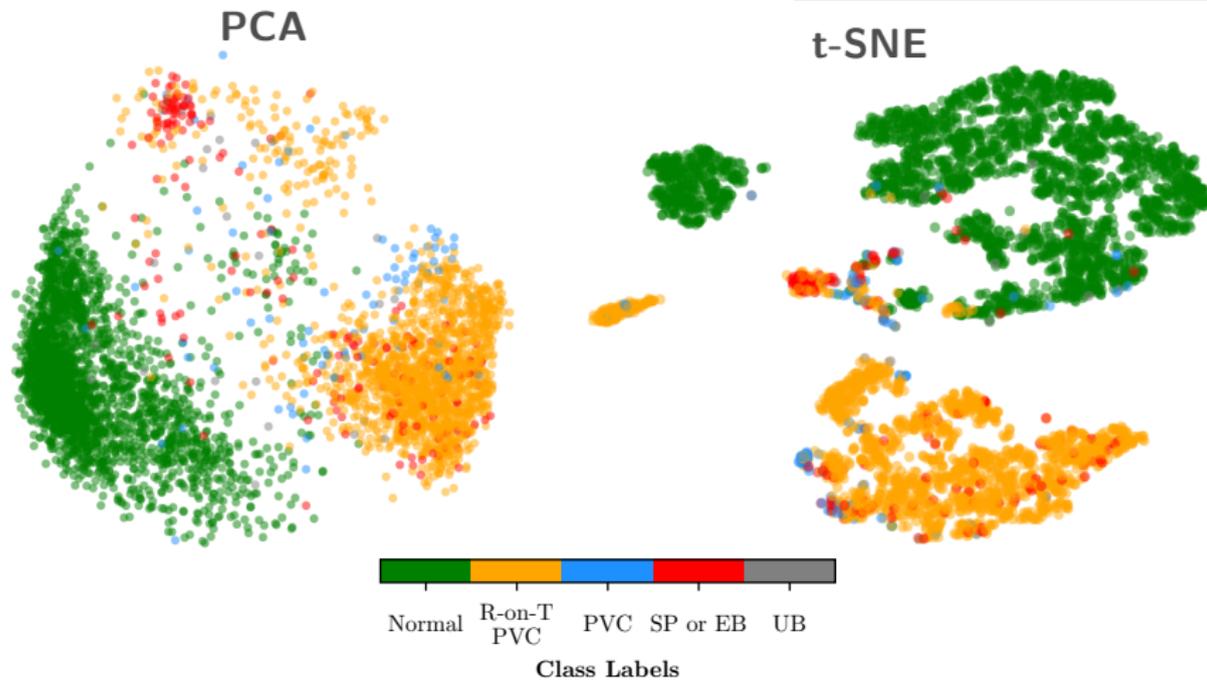


ECG5000



Each datapoint \rightarrow a sequence of length T

$T = 140$ (one heartbeat)
 $d_z = 5$



Metric	Hierarchical	Spectral	<i>k</i> -Means++	Wasserstein	SVM
AUC	0.9569	0.9591	0.9591	0.9819	0.9836
Accuracy	0.9554	0.9581	0.9596	0.9510	0.9843
Precision	0.9585	0.9470	0.9544	0.9469	0.9847
Recall	0.9463	0.9516	0.9538	0.9465	0.9843
F_1 -score	0.9465	0.9474	0.9522	0.9461	0.9844

Source	S/U	Model	AUC	Acc	F_1
Proposed	S	VRAE+SVM	0.9836	0.9843	0.9844
	U	VRAE+Clust/W	0.9819	0.9596	0.9522
Lei <i>et al.</i> , 2017	S	SPIRAL-XGB	0.9100	-	-
Karim <i>et al.</i> , 2017	S	F-t ALSTM-FCN	-	0.9496	-
Malhotra <i>et al.</i> , 2017	S	SAE-C	-	0.9340	-
Liu <i>et al.</i> , 2018	U	oFCMdd	-	-	0.8084

Unsupervised

Supervised

Metric	Hierarchical	Spectral	<i>k</i> -Means++	Wasserstein	SVM
AUC	0.9569	0.9591	0.9591	0.9819	0.9836
Accuracy	0.9554	0.9581	0.9596	0.9510	0.9843
Precision	0.9585	0.9470	0.9544	0.9469	0.9847
Recall	0.9463	0.9516	0.9538	0.9465	0.9843
F_1 -score	0.9465	0.9474	0.9522	0.9461	0.9844

Source	S/U	Model	AUC	Acc	F_1
Proposed	S	VRAE+SVM	0.9836	0.9843	0.9844
	U	VRAE+Clust/W	0.9819	0.9596	0.9522
Lei <i>et al.</i> , 2017	S	SPIRAL-XGB	0.9100	-	-
Karim <i>et al.</i> , 2017	S	F-t ALSTM-FCN	-	0.9496	-
Malhotra <i>et al.</i> , 2017	S	SAE-C	-	0.9340	-
Liu <i>et al.</i> , 2018	U	oFCMdd	-	-	0.8084

Unsupervised

Supervised

Metric	Hierarchical	Spectral	<i>k</i> -Means++	Wasserstein	SVM
AUC	0.9569	0.9591	0.9591	0.9819	0.9836
Accuracy	0.9554	0.9581	0.9596	0.9510	0.9843
Precision	0.9585	0.9470	0.9544	0.9469	0.9847
Recall	0.9463	0.9516	0.9538	0.9465	0.9843
F_1 -score	0.9465	0.9474	0.9522	0.9461	0.9844

Source	S/U	Model	AUC	Acc	F_1
Proposed	S	VRAE+SVM	0.9836	0.9843	0.9844
	U	VRAE+Clust/W	0.9819	0.9596	0.9522
Lei <i>et al.</i> , 2017	S	SPIRAL-XGB	0.9100	-	-
Karim <i>et al.</i> , 2017	S	F-t ALSTM-FCN	-	0.9496	-
Malhotra <i>et al.</i> , 2017	S	SAE-C	-	0.9340	-
Liu <i>et al.</i> , 2018	U	oFCMdd	-	-	0.8084

The proposed approach is:

- ▶ Effective on detecting anomalies in time series data;
- ▶ Suitable for both **univariate and multivariate** data;
- ▶ **Efficient**: inference and anomaly scores computation is fast;
- ▶ Works with other kinds of sequential data (e.g., **text, videos**);
- ▶ Extensible to a **multi-class framework** that allows discrimination between anomalies.

Unsupervised Anomaly Detection in Energy Time Series Data using Variational Recurrent Autoencoders with Attention

João Pereira
 Author for System and Robotics, Instituto Superior Técnico
 University of Lisbon
 Lisbon, Portugal
 joao.p.pereira.pereira@tecnico.ulisboa.pt

Margarida Silveira
 Author for System and Robotics, Instituto Superior Técnico
 University of Lisbon
 Lisbon, Portugal
 msilveira@tecnico.ulisboa.pt

Abstract—In the age of big data, time series are being generated in massive amounts. In the energy field, smart grids are enabling a computational data acquisition with the integration of sensors and smart devices. In the context of renewable energies, there has been an increasing interest in solar generation energy generation. These installations are often integrated with smart sensors that measure the energy production. Such amount of data collected makes the quest for developing smart monitoring systems that can detect anomalous behavior in these systems, trigger alerts and enable maintenance operations.

In this paper, we propose a generic, unsupervised and scalable framework for anomaly detection in time series data, based on a variational recurrent autoencoder. Furthermore, we introduce attention in the model, to account of a variational self-attention mechanism (VAM), to improve the performance of the encoding/decoding process. Afterwards, we perform anomaly detection based on the probabilistic reconstruction errors provided by our model.

Our results on solar energy generation time series show the ability of the proposed approach to detect anomalous behavior in time series data, while providing contextual and interpretable representations. Since it does not need labels to be trained, our methodology enables the application for anomaly detection in energy time series data and beyond.

Index Terms—Anomaly Detection, Variational Recurrent Autoencoder, Attention, Solar Photovoltaic Energy

I. INTRODUCTION

One of the key assets of the smart grid is the data it collects,

too expensive, while it requires domain knowledge from experts in the field. The lack of labels is, indeed, one of the reasons why anomaly detection has been such a great challenge for researchers and practitioners.

Furthermore, some of the proposed methods do not consider the sequential nature of the data by assuming it is independent in time. Smart grid data is often segmented by sensor and monthly time series and, hence, it is crucial to take into account the order and structure of the data.

The main contributions of this work can be summarized as follows:

- Unsupervised reconstruction-based model using a variational autoencoder with attention module and decoder;
- Variational self-attention mechanism to improve the encoding/decoding process;
- Generic framework for anomaly detection in time series data;
- Application to solar photovoltaic generation time series.

II. BACKGROUND

In this section, we review autoencoder, recurrent neural networks, attention mechanisms and autoencoder-based anomaly detection.

Learning Representations from Healthcare Time Series Data for Unsupervised Anomaly Detection

João Pereira
 Author for System and Robotics
 University of Lisbon
 Lisbon, Portugal
 joao.p.pereira.pereira@tecnico.ulisboa.pt

Margarida Silveira
 Author for System and Robotics, Instituto Superior Técnico
 University of Lisbon
 Lisbon, Portugal
 msilveira@tecnico.ulisboa.pt

Abstract—The amount of time series data generated in Healthcare is growing very fast and so is the need for methods that can analyze these data, detect anomalies and provide meaningful insights. However, most of the data available is unlabeled and, therefore, anomaly detection in this scenario has been a great challenge for researchers and practitioners.

Recently, unsupervised representation learning with deep generative models has been applied to find representations of data, without the need for big labeled datasets. Motivated by these success, we propose an unsupervised framework for anomaly detection in time series data. In our method, both representation learning and anomaly detection are fully unsupervised. In addition, the training data may contain anomalous data. We first learn representations of this series using a Variational Recurrent Autoencoder. Afterwards, based on these representations, we detect anomalous time series using Clustering and the Wasserstein Distance.

Our results on the publicly available MIMIC-III electrocardiogram dataset show the ability of the proposed approach to detect anomalous behavior in a fully unsupervised fashion, while providing structured and interpretable data representations. Furthermore, our approach outperforms previous supervised and unsupervised methods on this dataset.

Index Terms—Variational Recurrent Autoencoder, Representation Learning, Clustering, Electrocardiogram.

I. INTRODUCTION

In this work, we propose an unsupervised framework for anomaly detection in sequential data, based on representation learning using a Variational Recurrent Autoencoder and anomaly detection in the representation space via Clustering and the Wasserstein distance [1].

This paper is organized as follows. We start by reviewing Autoencoders, Variational Autoencoders and Recurrent Neural Networks. Then, we present a summary of recent approaches to anomaly detection in time series data. Afterwards, we introduce our proposed representation learning model and detection methodology. Finally, we present and analyze the results obtained with our model in electrocardiogram (ECG) time series.

Our contributions in this work can be summarized as:

- Unsupervised representation learning of time series data through a Variational Recurrent Autoencoder;
- Latent space-based detection using Clustering and the Wasserstein distance.

II. BACKGROUND

In this section, we review Autoencoders, Variational Autoencoders and Recurrent Neural Networks, including Long Short-Term Memory Networks.

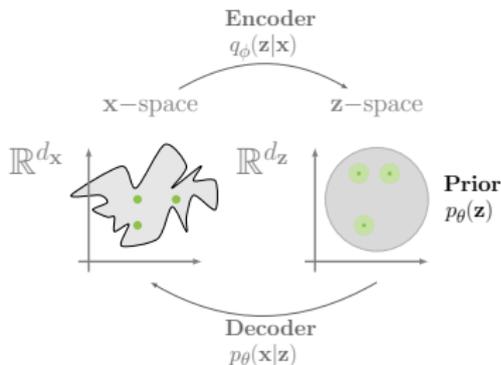
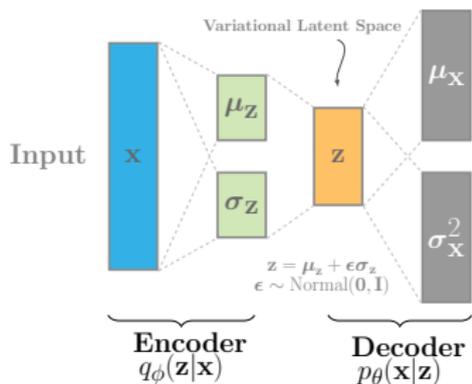
Accepted for Oral Presentation at the
 17th IEEE International Conference on
 Machine Learning and Applications

Submitted to the
 IEEE International Conference on Big
 Data and Smart Computing

Thank you for your attention!

The Variational Autoencoder

- ▶ Deep generative model rooted in Bayesian inference



$$p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z}$$

$$p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{x}|\mathbf{z})}{p_{\theta}(\mathbf{x})}$$

- ▶ **Objective:** Maximize the Evidence Lower Bound (ELBO)

$$\log p_{\theta}(\mathbf{x}) \geq \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{:= \mathcal{L}_{\text{ELBO}}(\theta, \phi; \mathbf{x})} - \mathcal{D}_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z}))$$

Kingma & Welling, **Auto-Encoding Variational Bayes**, ICLR'14

Rezende *et al.*, **Stochastic Backpropagation and Approximate Inference in Deep Generative Models**, ICML'14

Derivation of the Evidence Lower Bound (ELBO)

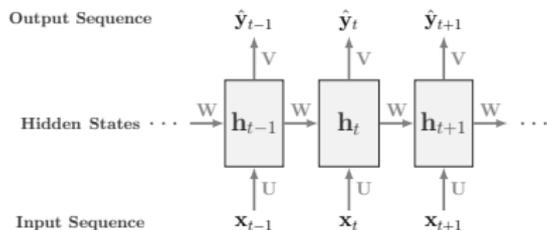
$$\begin{aligned}\log p_\theta(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x})] \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p_\theta(\mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p_\theta(\mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})} \right] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} \right] \\ &= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z}))}_{=\mathcal{L}_{\text{ELBO}}(\theta, \phi; \mathbf{x})} + \underbrace{\mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z}|\mathbf{x}))}_{\geq 0}\end{aligned}$$

What if data are not i.i.d. in time?

(e.g., time series, text, videos)

RNNs capture the temporal dependencies of the data

- ▶ Real-valued hidden state \mathbf{h}_t
- ▶ Feedback connection
- ▶ Parameters shared across timesteps



$$\mathbf{h}_t = f(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{h}_{t-1})$$

f is typically a tanh or sigmoid

Long Short-Term Memory Network

- ▶ Proposed to solve the vanishing gradient problem
- ▶ New cell and three gates

- ▶ Updates:

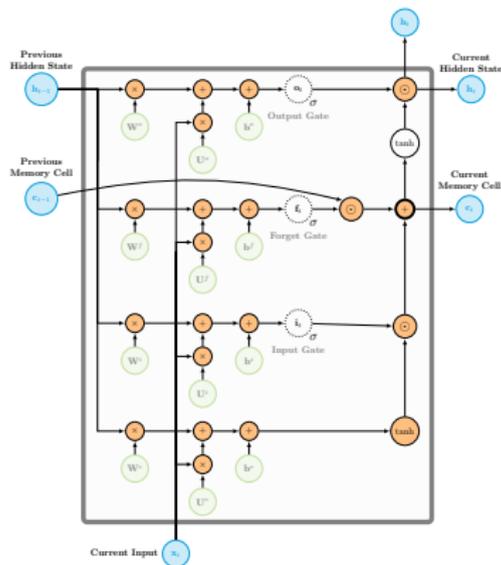
$$\mathbf{i}_t = \sigma(\mathbf{W}_i \mathbf{h}_{t-1} + \mathbf{U}_i \mathbf{x}_t + \mathbf{b}_i)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_f \mathbf{h}_{t-1} + \mathbf{U}_f \mathbf{x}_t + \mathbf{b}_f)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_o \mathbf{h}_{t-1} + \mathbf{U}_o \mathbf{x}_t + \mathbf{b}_o)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{W}_c \mathbf{h}_{t-1} + \mathbf{U}_c \mathbf{x}_t + \mathbf{b}_c)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

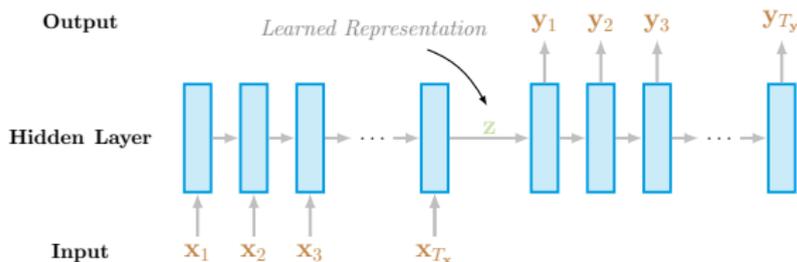


Hochreiter *et al.*, **Long-Short Term Memory**, Neural Computation'97

Graves *et al.*, **Bidirectional LSTM Networks for Improved Phoneme Classification and Recognition**, ICANN'05

Sequence to Sequence (Seq2Seq) models

Map sequences to sequences is useful.



Sutskever *et al.*, **Sequence to Sequence Learning with Neural Networks**, NIPS'14

Cho *et al.*, **Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation**, NIPS'14

Classification Metrics

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F_1\text{-score} = 2 \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

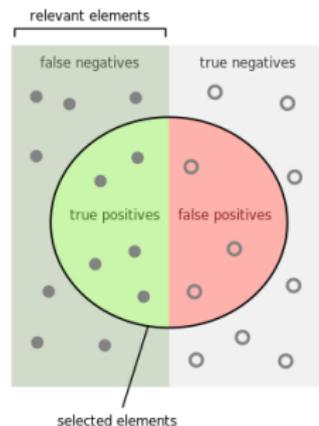
Legend:

TP : True Positives

TN : True Negatives

FP : False Positives

FN : False Negatives



How many selected items are relevant?

$$\text{Precision} = \frac{\text{green}}{\text{green} + \text{red}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{green}}{\text{green}}$$

Classification Metrics

AUC: Area under the Receiver Operating Characteristic Curve
(average precision over recalls)

