

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

João Pereira & Margarida Silveira

Signal and Image Processing Group
Institute for Systems and Robotics
Instituto Superior Técnico (Lisbon, Portugal)



Orlando, Florida, USA
December 19th, 2018

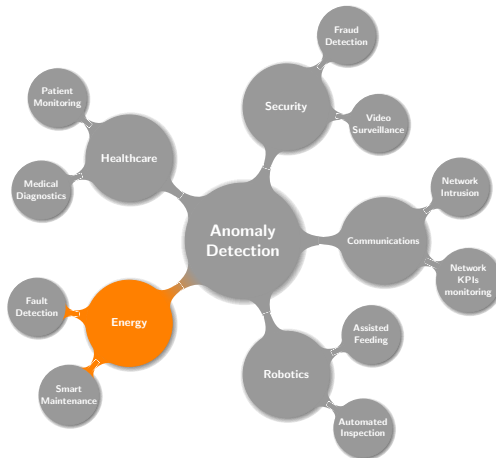


Introduction

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



Introduction



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

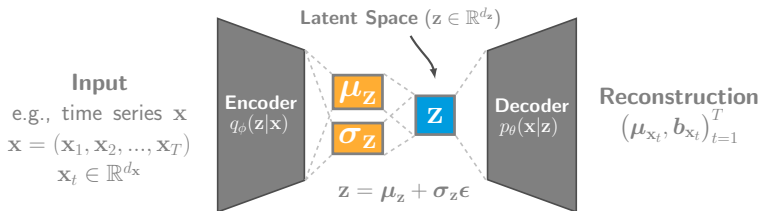


- ▶ Time series have temporal structure/dependencies

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

The Principle in a Nutshell

- Based on a **Variational Autoencoder**¹²

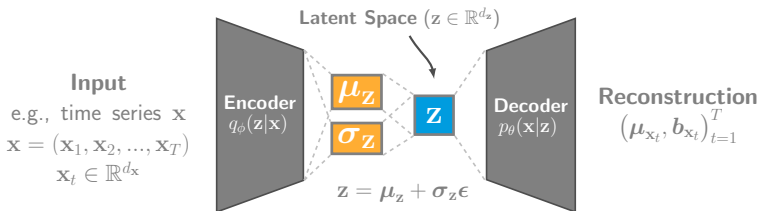


¹Kingma & Welling, **Auto-Encoding Variational Bayes**, ICLR'14

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

The Principle in a Nutshell

- ▶ Based on a **Variational Autoencoder**¹²



- ▶ Train a VAE on data with mostly **normal** patterns;
- ▶ It reconstructs well **normal** data, while it fails to reconstruct **anomalous** data;
- ▶ The **quality of the reconstructions** is used as anomaly score.

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

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

Proposed Approach

Reconstruction
Model

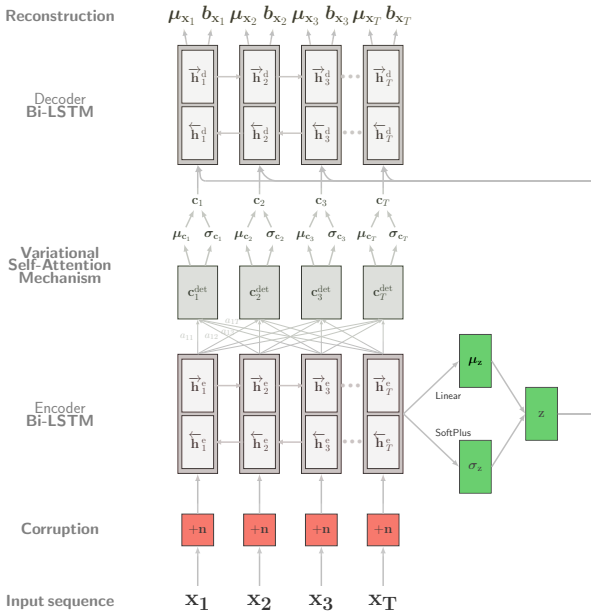
Detection

Proposed Approach

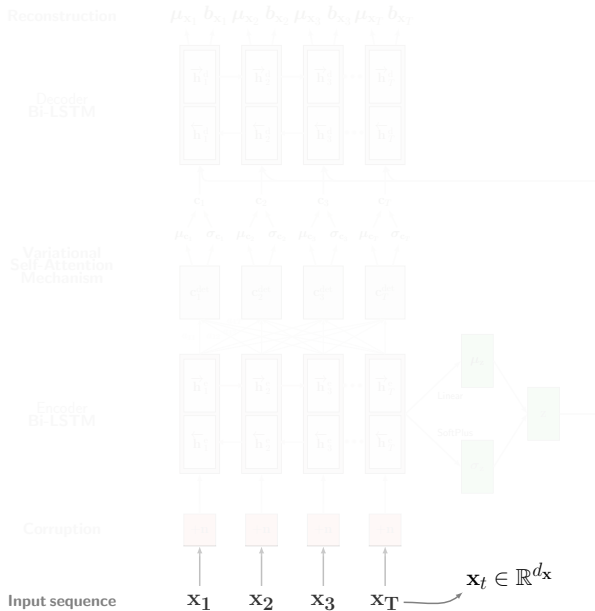
**Reconstruction
Model**

Detection

Reconstruction Model



Reconstruction Model



Reconstruction Model

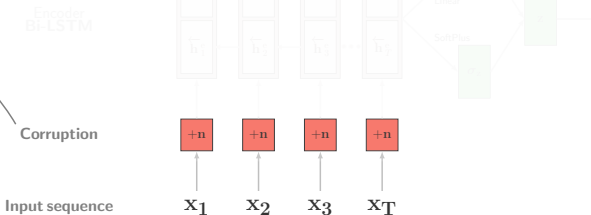
Denosing Autoencoding Criterion

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



Reconstruction Model

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

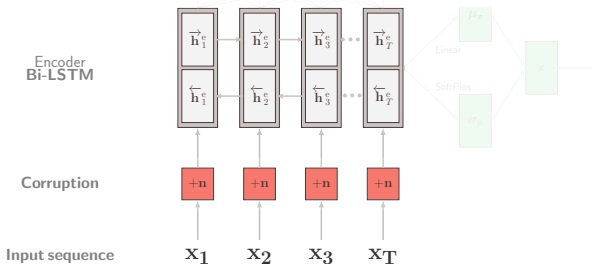
Bidirectional Long-Short Term Memory network

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

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

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

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



Reconstruction Model

Variational Latent Space

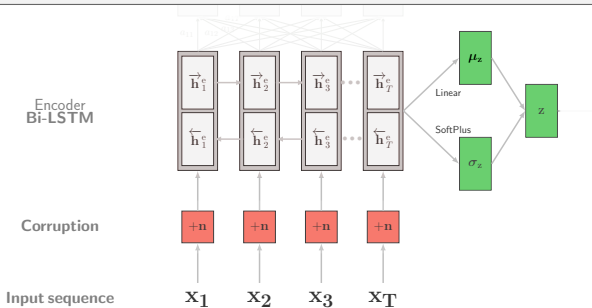
Variational parameters derived using neural networks

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

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'14



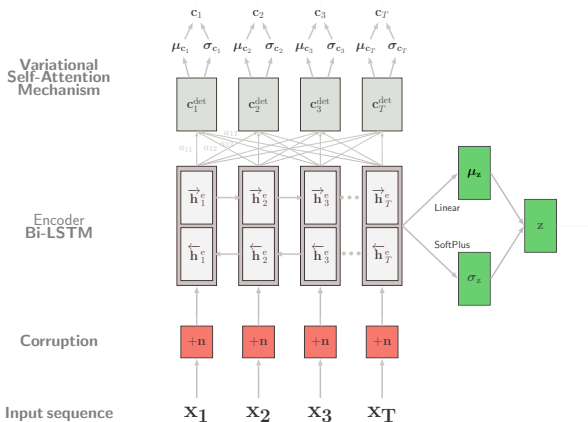
Reconstruction Model

Variational Self-Attention

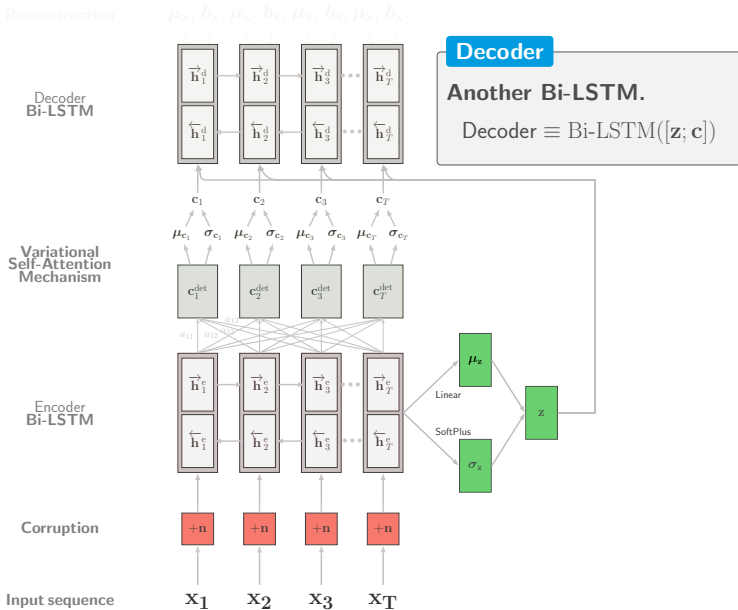
Combines self-attention with variational inference.

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

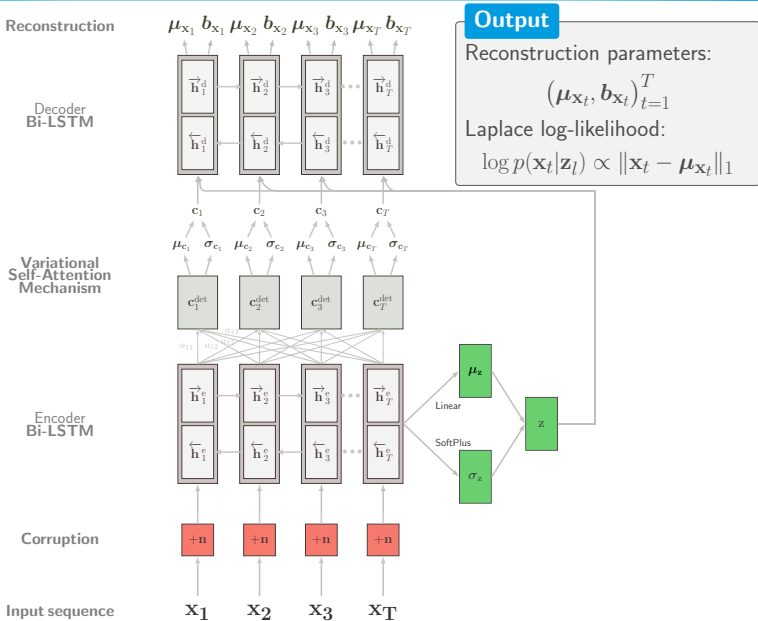
Vaswani *et al.*, Attention is All You Need, NIPS'17



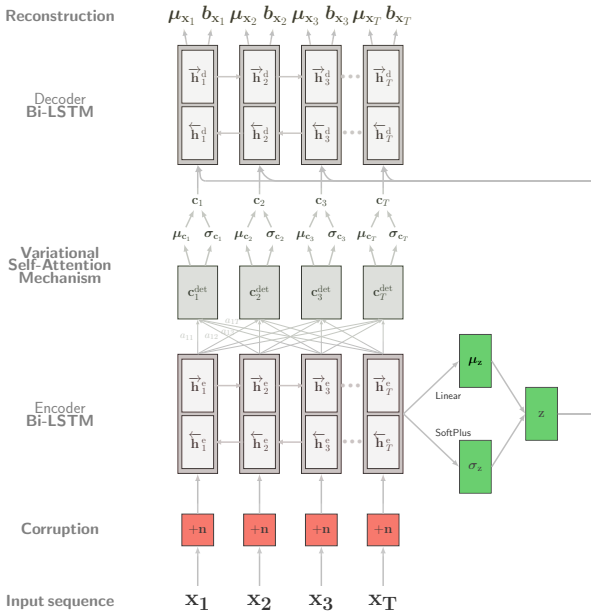
Reconstruction Model



Reconstruction Model



Reconstruction Model



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

Loss Function

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

\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

Reconstruction
Model

Detection

Proposed Approach

Reconstruction
Model

Detection

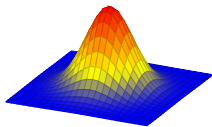
Anomaly Scores

Anomaly Scores

- Use the **reconstruction error** as anomaly score.

$$\text{Anomaly Score} = \mathbb{E}_{\mathbf{z}_l \sim q_\phi(\mathbf{z}|\mathbf{x})} \left[\left\| \mathbf{x} - \underbrace{\mathbb{E}[p_\theta(\mathbf{x}|\mathbf{z}_l)]}_{\boldsymbol{\mu}_x} \right\|_1 \right]$$

$$q_\phi(\mathbf{z}|\mathbf{x}) = \text{Normal}(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z^2 \mathbf{I})$$



$$\mathbf{z}_l \sim q_\phi(\mathbf{z}|\mathbf{x})$$

Anomaly Scores

- ▶ Use the **reconstruction error** as anomaly score.

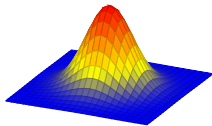
$$\text{Anomaly Score} = \mathbb{E}_{\mathbf{z}_l \sim q_\phi(\mathbf{z}|\mathbf{x})} \left[\left\| \mathbf{x} - \underbrace{\mathbb{E}[p_\theta(\mathbf{x}|\mathbf{z}_l)]}_{\boldsymbol{\mu}_x} \right\|_1 \right]$$

- ▶ Take the **variability of the reconstructions** into account.

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

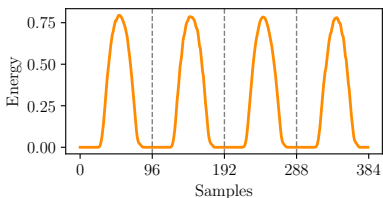
$$q_\phi(\mathbf{z}|\mathbf{x}) = \text{Normal}(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z^2 \mathbf{I})$$

$$\mathbf{z}_l \sim q_\phi(\mathbf{z}|\mathbf{x})$$



Experiments & Results

Solar PV Generation



(Production in a day without clouds)



- ▶ Provided by

c|side

- ▶ Recorded every 15min (96 samples per day)
- ▶ Data normalized to the installed capacity
- ▶ Daily seasonality
- ▶ Unlabelled

Variational Latent Space

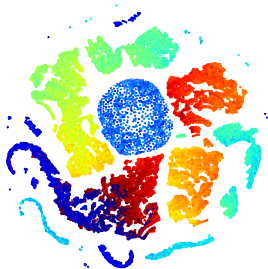
z-space in 2D ($\mathcal{X}_{\text{train}}^{\text{normal}}$)

$T = 32$ (< 96)

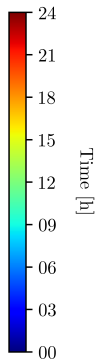
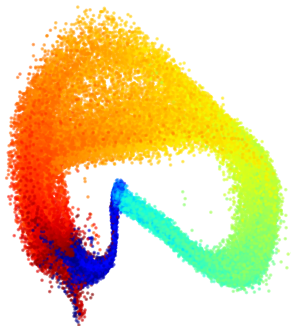
$d_z = 3$

online mode

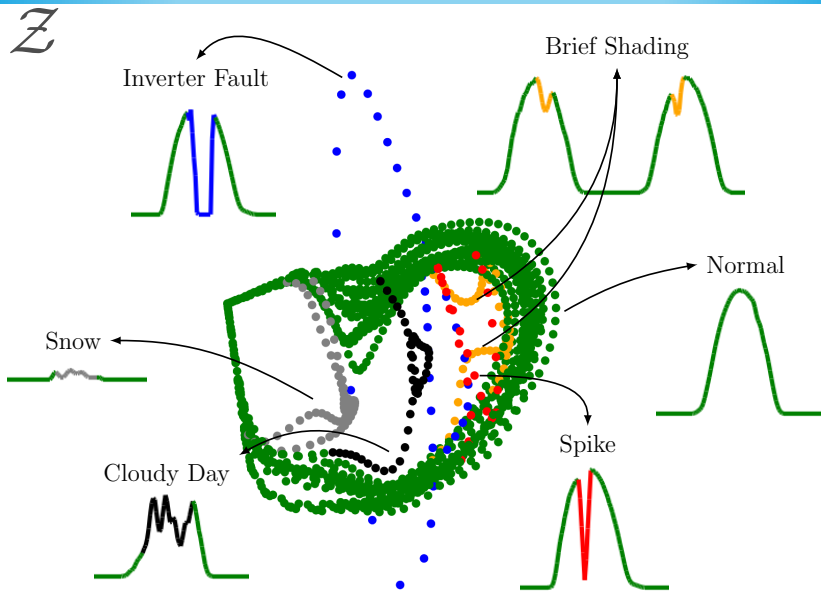
t-SNE



PCA



Finding Anomalies in the Latent Space



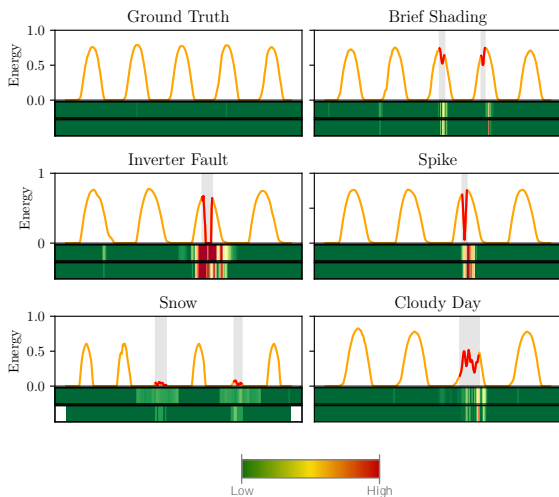
Anomaly Scores

Reconstruction Error

(top bar)

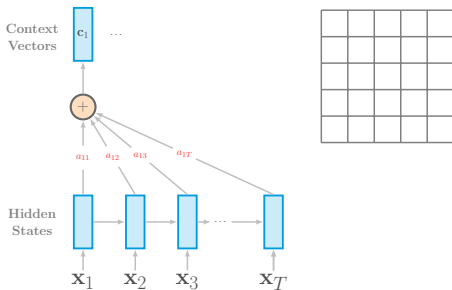
–Reconstruction Probability

(bottom bar)



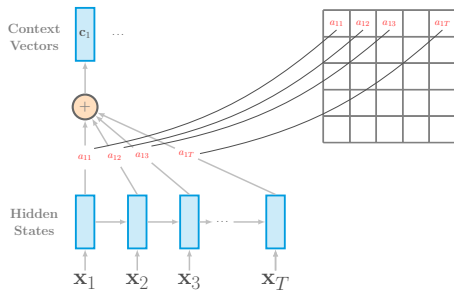
Attention Visualization

Attention Map

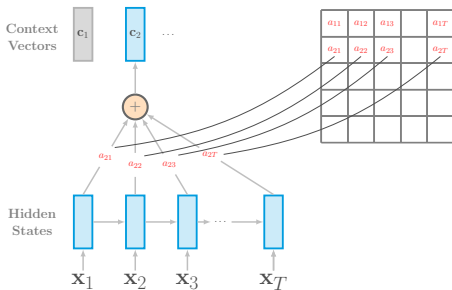


Attention Visualization

Attention Map

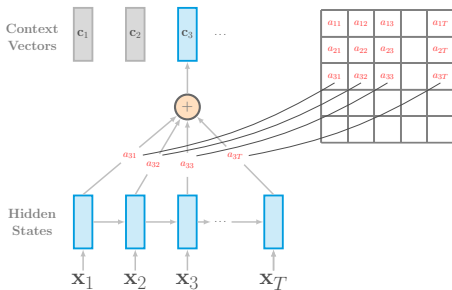


Attention Map

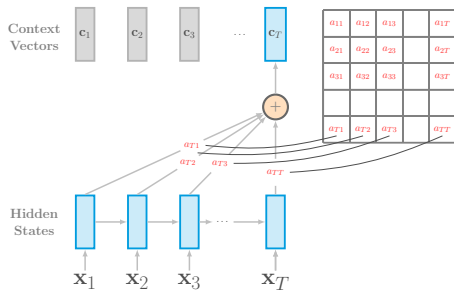


Attention Visualization

Attention Map



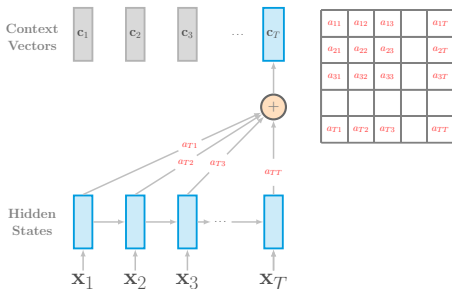
Attention Map



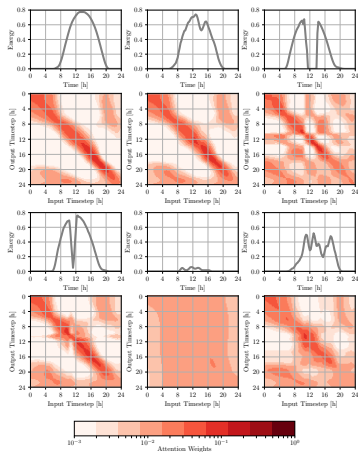
Examples

Attention Visualization

Attention Map



Examples



Conclusions & Future Work

- ▶ Effective on detecting anomalies in time series data;
- ▶ Unsupervised;
- ▶ Suitable for both **univariate and multivariate** data;
- ▶ **Efficient**: inference and anomaly scores computation is fast;
- ▶ **General**: works with other kinds of sequential data (e.g., **text, videos**);
- ▶ Exploit the usefulness of the **attention maps** for detection;
- ▶ Make it robust to changes of the normal pattern over time.

Acknowledgements



Thank you for your attention!

mail@joao-pereira.pt
www.joao-pereira.pt

