

MGADN: A Multi-task Graph Anomaly Detection Network for Multivariate Time Series

Weixuan Xiong

Yun Qin

Xiaochen Sun

Tianjin University, Tianjin, China

WEIXUAN@TJU.EDU.CN

QINYUNNN@163.COM

XCSUN@TJU.EDU.CN

Abstract

Anomaly detection of time series, especially multivariate time series (time series with multiple sensors), has been focused on for several years. Though existing method has achieved great progress, there are several challenging problems to be solved. Firstly, existing method including neural network only concentrate on the relationship in terms of timestamp. To be exact, they only want to know how does the data in the past influence which in the future. However, one sensor sometimes intervenes in other sensor such as the speed of wind may cause decrease of temperature. Secondly, there exist two categories of model for time series anomaly detection: prediction model and reconstruction model. Prediction model is adept at learning timely representation while short of capability when faced with sparse anomaly. Conversely, reconstruction model is opposite. Therefore, how can we efficiently get the relationship both in terms of both timestamp and sensors becomes our main topic. Our approach uses GAT, which is originated from graph neural network, to obtain connection between sensors. And LSTM is used to obtain relationships timely. Our approach is also designed to be double headed to calculate both prediction loss and reconstruction loss via VAE (Variational Auto-Encoder). In order to take advantage of two sorts of model, multi-task optimization algorithm is used in this model.

Keywords: Multivariate Time Series, Anomaly Detection, GNNs, Multi-task Learning

1. Introduction

A time series refers to a series of sequences arranged in the order of time occurrence. Time series are now widely used in speech analysis, noise removal, and stock market analysis. Multivariate time series is those time series having multiple channels. Through time series analysis, we can also find potential connections and numerical features implicit in these time series data, so that we can better understand various scientific theories, and we can also discover social phenomena from the data. Nature time series anomaly detection can help us locate anomalies in time, including systemic errors, disasters, and hardware and software damages. And use this to timely check and fill gaps to avoid the expansion of losses. Therefore, how to establish an accurate and efficient anomaly detection model is also an important topic in the field of time series analysis.

In real world, many systems product huge amount of inter-related multivariate time series data. For instance, in climate prediction system, different climate sensors receive different kinds of climate data such as temperature, moisture and wind speed etc. Apparently, those data collected by different sensors is related to each other which means one change can result in others' change. Most relationships between sensors are non-linear. For

example, the surge of wind speed or sudden change of wind direction may bring massive cold air which in turns cause temperature decrease. This pattern can't be learned only from previous temperature.

In the past, traditional methods such as ARIMA and machine learning methods such as SVM(Support Vector Machines) are used to detect anomalies. Nevertheless, with the rapid growth of dimensionality and complexity of data from more sensors, traditional time series analysis methods are not appropriate anymore. Thus, deep learning methods become more and more popular and practical because of their strong ability of fitting different kinds of data. In terms of time series anomaly detection, two major deep learning methods are prediction model and reconstruction model. Fundamentally, AE(Auto Encoder) is a fairly practical model which use reconstruction error to discover the occurrence of anomaly. Its upgraded version-VAE(Variable Auto Encoder) regulate the distribution of latent space and use Kullback-Leibler divergence and binary cross entropy to calculate loss function which obviously strengthen its reconstruction power. More recently, Generative Adversarial Networks (GANs), which is also acted as a reconstruction model, is used in the area. Construction models prove their competence to detect anomalies especially when anomalous labels are sparse, however proves their weakness when excavating timely relationships. On the contrary, prediction models such as RNN or LSTM-based approaches take full advantage of their sequential input pattern to learn timely relationships. But when faced with sparse anomaly, the accuracy of prediction models decrease. In addition both two models cannot explicitly learn the relationships between sensors. This limits their power to detect and explain deviations from such relationships when anomalous events happen.

In order to better learn the connections between sensors and make it contribute to the power of anomaly detection, Graph Neural Networks(GNNs) [Defferrard et al. \(2016\)](#) should be considered. GNNs, including graph convolution networks (GCNs) [Kipf and Welling \(2016\)](#), graph attention networks (GATs) [Veličković et al. \(2017\)](#), have shown its great ability to learn relationships between in graph structured data. However, most GNNs take data with pre-defined graph as input which is not applicable in multivariate time series data. Therefore, how to utilize GNNs in dealing with multi-time series without any prior information becomes our major topic. The first approach that fits GNNs on time series analysis is MTGNN proposed by [Wu et al. \(2020\)](#). While MTGNN is for time series prediction, several GNNs approaches came up for not only prediction, but for anomaly detection or classification. However, they are either prediction model or reconstruction model, in other words, one headed model.

Thus, in this paper, we proposed our MGADN: A Multi-task Graph Anomaly Detection Network, which combines both prediction model and reconstruction model and uses GNNs to discover inherent relationship between sensors. MGADN have four main components:(1).**Shared Weight Layer**: provide shared parameters of the whole model and acquire timely pattern through certain designed structure.(2).**Construction head**: using VAE to reconstruct input and calculating reconstruction error. (3).**Prediction head**: using GNNs to learn graph from multi-time series and slide window methods to make prediction and calculating prediction error(deviation between ground truth and predicted value).(4).**Multi-task Optimization and Deviation Scoring**: calculate loss from two heads with multi-task optimization methods and identifies and explains deviations from the learned sensor relationships in the graph.

In conclusion, this work makes contributions below:

- Compared with the existing one-head model such as MTGNN [Wu et al. \(2020\)](#) and GDN [Deng and Hooi \(2021\)](#), we propose MGADN, a new two-headed and graph neural network based approach that exploit advantage of both prediction and reconstruction model and learn inter-relationship between sensors. The GNNs module can automatically get graph structure from dataset itself.
- MGADN is modeled as a multi-task optimization method. When dealing with prediction loss and reconstruction loss, this multi-task learning assignment is transferred into multi-task optimization problem by calculating Pareto Optimality. This is for solving conflict between prediction task and reconstruction task to a certain extent.
- We conduct experiments on three open datasets with ground truth anomalies. Our results demonstrate that MGADN detects anomalies more accurately than baseline approaches. Some other experiments such as ablation study and visualization work will prove efficiency of every individual module.

2. Related Work

In this section, we will introduce several model related to multivariate time series modelling. In addition, since our approach is combined by graph neural network and multi-task learning, methods about these two topics will also be summarized in this section.

Multivariate Time Series Modelling: The earliest time series analysis method is the auto-regressive (AR) model proposed by British statistician [Udny Yule \(1927\)](#). In 1931, mathematician [Walker \(1931\)](#) established the Moving Average Model (Moving Average Model, MA). In the 1970s, the famous American statistician [Box et al. \(2015\)](#), scientists who made outstanding contributions to time series analysis, proposed a method for abnormal detection of non-stationary time series. The non-stationary time series is transformed into a stationary series through appropriate order operations, and then the method of ARMA model is used. That is, the differential autoregressive moving average model (Autoregressive Integrated Moving Average Model, ARIMA). In 1980, [Lütkepohl \(2013\)](#) proposed the vector autoregressive model (Vector Autoregressive Models VAR). The VAR model is a generalization of the AR model. Since the AR model can only be applied to the scenario of univariate time series analysis, it has great limitations, and VAR extends the AR model to multivariate time series analysis.

Since 20s, the rapid rise of deep learning has brought more attention to neural network methods. [Zaremba et al. \(2014\)](#) proposed Neural Network (Recurrent Neural Network RNN), as a cyclic feedback neural network framework, is often used in the processing of sequence data. It can fully capture the correlation in time series across time stamps. And its revised versions such as Long Short Term Memory Networks(LSTMs) and Gated Recurrent Unit Networks(GRUs) are proved to be practical in multivariate time series analysis [Karim et al. \(2017\)](#); [Kumar et al. \(2018\)](#). Generative Adversarial Networks(GANs) [Li et al. \(2019\)](#); [Zhou et al. \(2019\)](#), Variable Auto-Encoder(VAE) and other reconstruction model are majorly used in anomaly detection.

Graph Neural Networks: In recent years, GNNs are considered as good methods to deal with graph structured data. The iterative mode of GNNs [Defferrard et al. \(2016\)](#)

is that the hidden state of one node is influenced by its neighbors and its own previous hidden state. Graph Convolution Networks(GCN) Kipf and Welling (2016) defines Fourier transformation on graph and upgrade its parameters by absorbing nodes' one-step neighbors' representation. Graph Attention Networks(GAT) Veličković et al. (2017) include the thesis of attention which calculate attention weight of each neighbors to get the importance rank of them during aggregation process. To move forward a single step, some works prove the ability of GNNs on modelling large-scale multivariate time series data. For instance, STGCN proposed in 2017 have succeeded in predicting traffic flows Yu et al. (2017).

With regard to anomaly detection, GDN, which is a prediction model, proposed by Deng and Hooi (2021) use GATs to capture relationship between sensors. Moreover, GNNs can also be fixed into reconstruction model like VAE or GAN. Li et al. (2021) proposed Stackvae-G which combines VAE and GNNs to detect anomaly in multi-time series.

Multi-task Learning: First we give the definition of Multi-task Learning(MTL) from Zhang and Yang (2018):

- **Definition 1 MTL:** Given m learning tasks $\{\mathcal{T}_i\}_{i=1}^m$ where all the tasks or a subset of them $i = 1$ are related but not identical, multi-task learning aims to help improve the learning of a model for T_i by using the knowledge contained in the m tasks.

MTL can be divided into two categories: **hard parameter sharing** and **soft parameter sharing** Ruder (2017). Hard parameter sharing method set hidden layers which have shared parameters between all tasks. In soft parameter sharing on the other hand, each task has its own model with its own parameters. MGADN will use hard parameter sharing method. On the other hand, to balance different tasks, there are different ways of training set. One may set model train different tasks together with the a loss that is the weighted sum of different losses from diversified tasks. Another way is to train different tasks one by one which means several times of loss backward. Sener and Koltun (2018) treated MTL as a multi-task optimization problem, and proved their methods practical.

3. Proposed Model

3.1. Problem Statement

In this work, our input of model is a multivariate time series dataset with N sensors and T time steps. We can denote input as $[h^1, h^2, \dots, h^T]$. With certain time stamp t , $h^t \in R^N$. With respect to prediction model, our task is to use slide window with its size d to predict future value. The goal projection is: $f : [h^{t+1}, h^{t+2}, \dots, h^{t+d}] \rightarrow \hat{h}^{t+d+1}$. As for reconstruction model, we can also use slide window with the same size d . The reconstruction goal will be $g : [h^{t+1}, h^{t+2}, \dots, h^{t+d}] \rightarrow [\hat{h}^{t+1}, \hat{h}^{t+2}, \dots, \hat{h}^{t+d}]$. Moreover, we will calculate deviation score through prediction and reconstruction. Specifically, the deviation between prediction or reconstruction and ground truth will be transferred into error score which act as the criterion for anomaly diagnoses. That is, with a certain time t , the output of the model will be binary value $score(t)$, i.e, $score(t) \in \{0, 1\}$, where $score(t) = 1$ means occurrence of anomaly on time t .

3.2. Overview of MGADN

Our proposed model MGADN is composed of four major modules: **Shared Weight Layer**, **Graph Embedding Module**, **VAE head** and **Graph Attention Forecasting head**. The outputs of two heads will be transferred into deviation score and judged. Furthermore, because of our two-headed design, this model can be regarded as a multi-task learning framework. Due to this specified design, MGDA algorithm will be utilized in order to balance the training of two heads. Fig. 1 shows the overview framework of our proposed model.

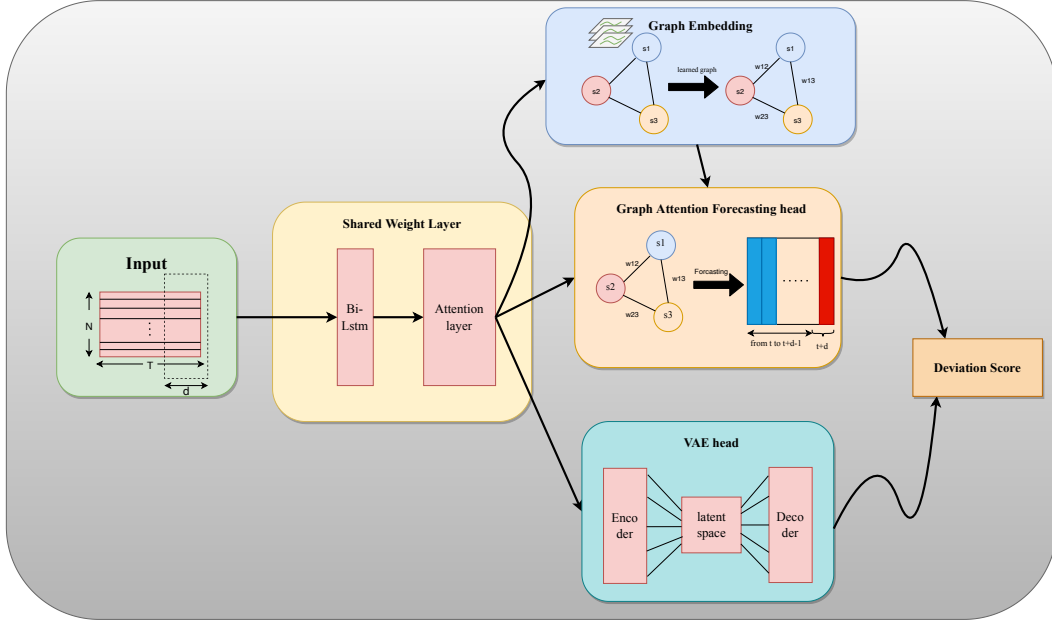


Figure 1: Overview of MGADN.

3.3. Shared Weighted Layer

The first module of our model is a layer whose weights are shared by the whole model. The shared layer is composed of a Bi-LSTM layer and a self attention layer. The framework of this layer is showed in fig2. Bi-LSTM acts as a time pattern collector, which can learn the relation within one slide window across time. When the number of hidden layer of Bi-LSTM is 1, the input $[h^{t-d+1}, h^{t-d+2}, \dots, h^t]$ will be transferred into:

$$\left[\frac{(s_{forward}^{t-d+1} + s_{backward}^{t-d+1})}{2}, \frac{(s_{forward}^{t-d+2} + s_{backward}^{t-d+2})}{2}, \dots, \frac{(s_{forward}^t + s_{backward}^t)}{2} \right],$$

Attention technique is widely used in natural language processing task. For instance, in translation task, self attention can take full advantage of forward and backward context by setting attention weight. Attention weight describes the importance of input. The main

step of attention technique is to find three variables which are Q(query), K(key), V(value). In our self attention layer, we assume $Q = K = V$.

Firstly, three weighting matrixes should be initialized as W^Q, W^K, W^V . Then the input which is also the output of Bi-LSTM: X will be operated below:

$$\begin{cases} Q = W^Q X \\ K = W^K X \\ V = W^V X \end{cases} \quad (1)$$

$$Out = V \times Softmax\left(\frac{QK^T}{\sqrt{d_k}}\right) \quad (2)$$

To clarify the design of shared weighted layer, Bi-LSTM layer capture forward and backward

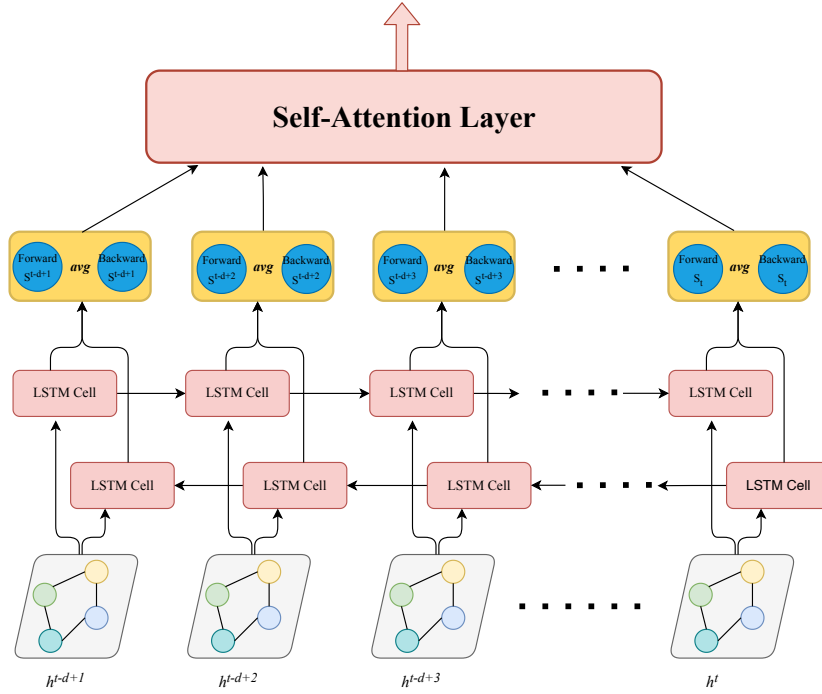


Figure 2: Shared Weighted Layer, which is composed of two parts:Bi-Lstm and Self Attention Layer.

timely pattern. Self attention in turns reinforce some time stamps by giving them higher weight. In fact these two sub-modules play similar roles in this model. However, one is the other one's enhancement. The architecture of shared weighted layer is showed in Fig. 2.

3.4. Graph Attention Forecasting Head

3.4.1. GRAPH STRUCTURE LEARNING

As is mentioned above, we introduce GNNs to learn relationships between sensors, which brings another problem: how can we transfer multi-time series into graph structured data? Wu et al. (2020) in 2020 firstly propose a common solution to this problem by considering every sensors as nodes with its hidden states $S_i \in R^T, i \in 1, 2, 3, \dots, N$ i.e N sensors with T time stamps. Moreover, there are several ways of forming graph structure, which is, in other words, calculating adjacency matrix in different methods. For example, we can design directed graph by setting adjacency matrix asymmetric or undirected graph in similar way. Here, our solution is to acquire symmetric adjacency matrix by calculating Pearson correlation coefficient.

Provided with N sensors, sensor i is potentially related to the rest of all sensors. We define the potential neighbors of sensor i as $P(i) = \{1, 2, \dots, i - 1, i + 1, \dots, N\}$. The output of shared weighted layer will be initially sent to an embedding layer. The outputs of embedding layer are denoted as v_i for $i = 1, 2, 3, \dots, N$. Calculations are as followed:

$$e_{ji} = \frac{\mathbf{v}_i^\top \mathbf{v}_j}{\|\mathbf{v}_i\| \cdot \|\mathbf{v}_j\|} \text{ for } j \in \mathcal{C}_i \quad (3)$$

$$A_{ji} = 1 \{j \in \text{TopK}(\{e_{ki} : k \in \mathcal{C}_i\})\} \quad (4)$$

That is to say, we firstly calculate the Pearson correlation coefficient between every two sensors. Secondly, for a certain sensor, we leave the most related k sensors as its neighbors.

3.4.2. GRAPH ATTENTION LAYER

In this section, we bring in a specified Graph Attention Layer proposed by Deng and Hooi (2021).

The input of this module is $\mathbf{x}^{(t)} := [h^{t-d}, h^{t-d+1}, \dots, h^{t-1}]$ when slide window is of d length. And we need to give the prediction of h^t . We denote that $x_i^{(t)} \in R^d$ is input feature of node i , $\mathcal{N}(i) = \{j \mid A_{ji} > 0\}$ is the neighborhood of node i according to adjacency A . We have graph attention mechanism:

$$\mathbf{g}_i^{(t)} = \mathbf{v}_i \oplus \mathbf{W}\mathbf{x}_i^{(t)} \quad (5)$$

$$\pi(i, j) = \text{LeakyReLU} \left(\mathbf{a}^\top \left(\mathbf{g}_i^{(t)} \oplus \mathbf{g}_j^{(t)} \right) \right) \quad (6)$$

$$\alpha_{i,j} = \frac{\exp(\pi(i, j))}{\sum_{k \in \mathcal{N}(i) \cup \{i\}} \exp(\pi(i, k))}, \quad (7)$$

$$\mathbf{z}_i^{(t)} = \text{ReLU} \left(\alpha_{i,i} \mathbf{W}\mathbf{x}_i^{(t)} + \sum_{j \in \mathcal{N}(i)} \alpha_{i,j} \mathbf{W}\mathbf{x}_j^{(t)} \right) \quad (8)$$

where $W \in R^{w \times d}$ is a trainable weight matrix which is also the embedding matrix before Graph Structure Learning layer. This weight matrix applies a learnable linear transformation towards every node. $\alpha_{i,j}$ is the attention weight between node i and node j . Then,

equation (8) provides iteration rule, which in turns let us obtain the representation of all nodes, namely $\{\mathbf{z}_1^{(t)}, \dots, \mathbf{z}_N^{(t)}\}$. The output of Graph Attention Layer can be described as:

$$\hat{\mathbf{h}}^{(t)} = f_{\theta} \left(\left[\mathbf{v}_1 \circ \mathbf{z}_1^{(t)}, \dots, \mathbf{v}_N \circ \mathbf{z}_N^{(t)} \right] \right) \quad (9)$$

where \circ means element-wise multiply, f_{θ} means stacked fully-connected layers with output dimension N .

3.5. VAE Head

VAE Head, which acts as a reconstruction model, is a Variable Auto Encoder module which is a updated version of AE(Auto Encoder). The underlying logic that AE can be competent for multi-time series anomaly detection is that compared with normal data, abnormal data is more difficult to be represented by low-dimensional latent layer, so when abnormal data is reconstructed, its reconstruction error will increase significantly, so as to determine the occurrence of abnormality. Furthermore, VAE encode input as a distribution while AE only encode input as a point in the latent space. This uncertain probability distribution provides more robust representation than which is provided by certain point. The whole process of VAE is as followed:

$$\mu, \sigma = W_{\mu} g_{\theta_1}(X), W_{\sigma} g_{\theta_1}(X) \quad (10)$$

$$Z = \mu + \sigma \times \epsilon \quad (11)$$

$$\hat{X} = f_{\theta_2}(Z) \quad (12)$$

Equation (10) is the process of encoder which transfer input data X into a Gaussian Distribution with its own expectation and variance. Equation (11) is to find latent variable from this distribution. Equation (12) is decoder process. The encoder and decoder processes are both realized by trainable linear transformation.

3.6. Multi-task Learning and Deviation Score

Multi-task Learning: At last, the outputs of two heads will be the output of prediction head: $\hat{\mathbf{h}}^{(t)}$, and the output of reconstruction head: $[\hat{\mathbf{r}}^{(t-d)}, \hat{\mathbf{r}}^{(t-d+1)}, \dots, \hat{\mathbf{r}}^{(t-1)}]$ with respect to the input: $[\mathbf{h}^{(t-d)}, \mathbf{h}^{(t-d+1)}, \dots, \mathbf{h}^{(t-1)}]$ and ground truth at time t : $\mathbf{h}^{(t)}$. The loss function is:

$$L_{\text{pred}} = \frac{1}{T_{\text{train}} - d} \sum_{t=d+1}^{T_{\text{train}}} \left\| \hat{\mathbf{h}}^{(t)} - \mathbf{h}^{(t)} \right\|_2^2 \quad (13)$$

$$\begin{aligned} L_{\text{recon}} &= KL(N(\mu, \sigma^2) \| N(0, I)) + \frac{1}{T_{\text{train}} - d} \sum_{t=d+1}^{T_{\text{train}}} \left| \hat{\mathbf{r}}^{(t)} - \mathbf{h}^{(t)} \right|_1 \\ &= \frac{1}{2} (-\log \sigma^2 + \mu^2 + \sigma^2 - 1) + \frac{1}{T_{\text{train}} - d} \sum_{t=d+1}^{T_{\text{train}}} \left| \hat{\mathbf{r}}^{(t)} - \mathbf{h}^{(t)} \right|_1 \end{aligned} \quad (14)$$

It should be pointed out that the original loss function of VAE will vibrate so much during training which causes the decrease of accuracy. Therefore, we replace Binary Cross

Entropy with L1 loss, which in turns solve this problem. Obviously, the model has become a multi-task learning framework(MTL). To take full advantage of this MTL structure, MGDA(Multiple Gradient Descent Algorithm) proposed by [Sener and Koltun \(2018\)](#) can be used for reference.

When faced with multiple distinct loss functions, there is an intuitionistic way. That is, assign different weights to each loss function. Given M different tasks, the overall loss function will be:

$$\min_{\theta^{sh}, \theta^1, \dots, \theta^M} \sum_{t=1}^M c^m \hat{\mathcal{L}}^m(\theta^{sh}, \theta^m) \quad (15)$$

where θ^{sh} is parameters shared by all tasks while θ^m is the task specific parameters. Consider two sets of solutions θ and $\bar{\theta}$ such that $\hat{\mathcal{L}}^{m_1}(\theta^{sh}, \theta^{m_1}) < \hat{\mathcal{L}}^{m_1}(\bar{\theta}^{sh}, \bar{\theta}^{m_1})$ and $\hat{\mathcal{L}}^{m_2}(\theta^{sh}, \theta^{m_2}) > \hat{\mathcal{L}}^{m_2}(\bar{\theta}^{sh}, \bar{\theta}^{m_2})$ for some tasks m_1 and m_2 . In other words, solution θ is better for task m_1 whereas $\bar{\theta}$ is better for m_2 . It is not possible to compare these two solutions without a pairwise importance of tasks, which is typically not available. To solve this problem, MGDA provides us with a alternative solution by calculating Pareto solution. Firstly, we can give the definition of Pareto optimality for MTL.

The optimization goal becomes:

$$\min_{\theta^{sh}, \theta^1, \dots, \theta^M} \mathbf{L}(\theta^{sh}, \theta^1, \dots, \theta^M) = \min_{\theta^{sh}, \theta^1, \dots, \theta^M} \left(\hat{\mathcal{L}}^1(\theta^{sh}, \theta^1), \dots, \hat{\mathcal{L}}^M(\theta^{sh}, \theta^M) \right)^\top \quad (16)$$

Definition 2: Pareto optimality for MTL

- A solution θ dominates $\bar{\theta}$ if $\hat{\mathcal{L}}^m(\theta^{sh}, \theta^m) \leq \hat{\mathcal{L}}^m(\bar{\theta}^{sh}, \bar{\theta}^m)$ for all task m and satisfying $\mathbf{L}(\theta^{sh}, \theta^1, \dots, \theta^M) \neq \mathbf{L}(\bar{\theta}^{sh}, \bar{\theta}^1, \dots, \bar{\theta}^M)$.
- A solution $\tilde{\theta}$ is called Pareto optimal if there is no other solutions dominate $\tilde{\theta}$.

Considering the Karush-Kuhn-Tucker conditions for both shared and task specific parameters, the optimization task becomes:

$$\min_{\alpha^1, \dots, \alpha^M} \left\{ \left\| \sum_{m=1}^M \alpha^m \nabla_{\theta^{sh}} \hat{\mathcal{L}}^m(\theta^{sh}, \theta^m) \right\|_2^2 \mid \sum_{m=1}^M \alpha^m = 1, \alpha^m \geq 0 \quad \forall m \right\} \quad (17)$$

[Désidéri \(2012\)](#) proved that either the solution to this optimization problem is 0 and the resulting point satisfies the KKT conditions, or the solution gives a descent direction that improves all tasks. Hence, the optimization problem becomes solving equation (17). [Sener and Koltun \(2018\)](#) also come up with an approximation version of MGDA, which is MGDA-UB by clarifying the inequation below:

$$\left\| \sum_{m=1}^M \alpha^m \nabla_{\theta^{sh}} \hat{\mathcal{L}}^m(\theta^{sh}, \theta^m) \right\|_2^2 \leq \left\| \frac{\partial \mathbf{Z}}{\partial \theta^{sh}} \right\|_2^2 \left\| \sum_{m=1}^M \alpha^m \nabla_{\mathbf{Z}} \hat{\mathcal{L}}^m(\theta^{sh}, \theta^m) \right\|_2^2 \quad (18)$$

where Z is the output of shared weighted layer. Then, the approximation of MGDA becomes:

$$\min_{\alpha^1, \dots, \alpha^M} \left\{ \left\| \sum_{m=1}^M \alpha^m \nabla_{\mathbf{Z}} \hat{\mathcal{L}}^m(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^m) \right\|_2^2 \mid \sum_{m=1}^M \alpha^m = 1, \alpha^m \geq 0 \quad \forall m \right\} \quad (19)$$

When it comes to our model, the optimization problem is:

$$\min_{\alpha \in [0,1]} \left\| \alpha \nabla_{\boldsymbol{\theta}^{sh}} \hat{\mathcal{L}}_{pred}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^1) + (1 - \alpha) \nabla_{\boldsymbol{\theta}^{sh}} \hat{\mathcal{L}}_{recon}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^2) \right\|_2^2 \quad (20)$$

Then, according to MGDA-UB, equation (20) becomes:

$$\min_{\alpha \in [0,1]} \left\| \alpha \nabla_{\mathbf{Z}} \hat{\mathcal{L}}_{pred}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^1) + (1 - \alpha) \nabla_{\mathbf{Z}} \hat{\mathcal{L}}_{recon}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^2) \right\|_2^2 \quad (21)$$

Equation (21) is a dimensional quadratic function of α which has an analytical solution:

$$a = \left[\frac{\left(\nabla_{\mathbf{Z}} \hat{\mathcal{L}}_{recon}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}_{recon}) - \nabla_{\mathbf{Z}} \hat{\mathcal{L}}_{pred}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}_{pred}) \right)^\top \nabla_{\mathbf{Z}} \hat{\mathcal{L}}_{recon}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}_{recon})}{\left\| \nabla_{\mathbf{Z}} \hat{\mathcal{L}}_{pred}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}_{pred}) - \nabla_{\mathbf{Z}} \hat{\mathcal{L}}_{recon}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}_{recon}) \right\|_2^2} \right] \quad (22)$$

$$\hat{\alpha} = \max(\min(a, 1), 0) \quad (23)$$

$$Loss = \hat{\alpha} \mathcal{L}_{pred} + (1 - \hat{\alpha}) \mathcal{L}_{recon} \quad (24)$$

Deviation Score: After acquiring the output of two heads, the deviation score calculation will be:

$$Err_i^{pred}(t) = \left| \mathbf{h}_i^{(t)} - \hat{\mathbf{h}}_i^{(t)} \right| \quad (25)$$

$$Err_i^{recon}(t) = \left| \mathbf{h}_i^{(t)} - \hat{\mathbf{r}}_i^{(t)} \right| \quad (26)$$

$$a_i^{pred}(t) = \frac{Err_i^{pred}(t) - \tilde{\mu}_i^{pred}}{\tilde{\sigma}_i^{pred}} \quad (27)$$

$$a_i^{recon}(t) = \frac{Err_i^{recon}(t) - \tilde{\mu}_i^{recon}}{\tilde{\sigma}_i^{recon}} \quad (28)$$

$$A(t) = \max(\max_i a_i^{pred}(t), \max_i a_i^{recon}(t)) \quad (29)$$

where $\tilde{\mu}_i, \tilde{\sigma}_i$ are the median and inter-quartile range across time stamps with respect to $Err_i(t)$. $A(t)$ is calculated as the maximum between prediction score and reconstruction score, by which model can retain its robustness and becomes more sensitive towards anomalies.

At last, a time stamp t is labelled as anomaly when $A(t)$ exceed a certain threshold, which is set to be the max of $A(t)$ within the validation data.

Table 1: The Properties of Datasets

Datasets	Features	Train	Test	Anomalies
MSL	27	1564	2048	21.90%
SWaT	51	49678	44991	12.13%
WADI	127	78457	17280	5.83%

4. Experiment

In this section, we will proceed several experiments for MGADN on three open datasets. Moreover, in order to clarify interpretability of our model, ablation study and some visualization works will be done.

4.1. Datasets

We first introduce datasets on which we conduct our experiments:

- **MSL:** This data represents real spacecraft telemetry data and anomalies from the Curiosity Rover on Mars (MSL).
- **SWaT:** The Secure Water Treatment (SWaT) dataset is from a water treatment test platform coordinated by Singapore’s Public Utility Board [Mathur and Tippenhauer \(2016\)](#).
- **WADI:** Water Distribution (WADI) [Ahmed et al. \(2017\)](#) is a distribution system comprising a larger number of water distribution pipelines which is the extension of SWaT.

4.2. Baselines

- **PCA (Principal Component Analysis):** finds a low-dimensional projection that keeps most of the variance within the dataset. The deviation score is the reconstruction error of this projection [Shyu et al. \(2003\)](#).
- **KNN (K Nearest Neighbors):** uses each point’s distance to its kth nearest neighbor as an deviation score [Angiulli and Pizzuti \(2002\)](#).
- **AE (Autoencoders):** An encoder-decoder model which uses reconstruction error as its deviation score [Aggarwal](#).
- **DAGMM (Deep Autoencoding Gaussian Model):** combines Gaussian Mixture Model with AE. It uses reconstruction error as its deviation score [Zhou et al. \(2019\)](#).
- **LSTM-VAE:** replaces the encoder and decoder which is composed of linear layers in VAE with LSTM. It uses reconstruction error as its deviation score [Park et al. \(2018\)](#).
- **MAD-GAN:** A GAN-like model, which uses LSTM-RNN as discriminator along with a reconstruction based anomaly detection approach [Li et al. \(2019\)](#).
- **GDN:** A prediction model, which combines GAT with anomaly detection approach and calculate deviation score at every predicted time stamps [Deng and Hooi \(2021\)](#).

Table 2: Experiment Result of Precision(%), Recall(%) and F1-score.

Dataset	MSL			SWaT			WADI		
	Prec.	Rec.	F1.	Prec.	Rec.	F1.	Prec.	Rec.	F1.
PCA	24.53	16.34	0.10	24.92	21.63	0.23	39.53	5.63	0.10
KNN	10.96	10.96	0.08	7.83	7.83	0.08	7.76	7.75	0.08
AE	53.92	67.63	0.60	72.63	52.63	0.61	34.35	34.35	0.34
DAGMM	49.24	70.32	0.29	27.46	69.52	0.39	54.44	26.99	0.36
LSTM-VAE	69.82	79.93	0.37	96.24	59.91	0.74	87.79	14.45	0.25
MAD-GAN	70.67	82.33	0.76	98.97	63.74	0.77	41.44	33.92	0.37
GDN	79.08	99.50	0.88	99.35	68.12	0.81	97.50	40.19	0.57
MGADN-ALT	78.31	98.81	0.87	98.45	69.44	0.84	92.97	49.33	0.64
MGADN-MGDA	80.80	98.87	0.89	99.50	76.95	0.87	94.16	50.98	0.66

4.3. Setup

We program our model with Pytorch version 1.5.1 with Cuda 11.1. PyTorch Geometric Library version 1.5.0 Fey and Lenssen (2019). The model is trained on Intel(R) Xeon(R) CPU E5-2699 v3 @ 2.30GHz and RTX2080TI. The models are trained using the Adam optimizer with learning rate 0.001. The number of epochs are at most 50. And we choose the number of neighbors k as 5(15, 30) for MSL(SWaT, WADI). The length of slide window will be 5 for all datasets.

4.4. Metrics:

In this section, we will introduce the evaluation metrics that is used in experiments. We use precision(prec), recall(rec) and F1 score(F1) on the test data and its ground truth to measure the capability of models. $Prec = \frac{TP}{TP+FP}$, $Rec = \frac{TP}{TP+FN}$ and $F1 = \frac{2Prec \times Rec}{Prec+Rec}$, where TN, TP, FN and FP represent true negative, true positive, false negative and false positive.

4.5. Experiment Results

4.5.1. ACCURACY

In table. 2, the results of experiments show that MGADN outperforms baselines in terms of precision, recall and F1-score. We also use two different ways to train MGADN: with MGDA and without MGDA (i.e MGADN-MGDA and MGADN-ALT). To specify, we train two heads alternately in MGADN-ALT. The result shows that MGDA-UB algorithm do make contributions in the model.

4.5.2. ABLATION STUDY

In this section, we proceed ablation study on dataset SWaT to justify the necessity of each module. The result shows that without MGDA-UB to balance two losses from two heads, the precision will slightly decrease but recall will decrease more. Dropping VAE head will cause decrease of recall. Deleting predicted head will cause decrease of all three metrics. Table. 3 shows the results of ablation study.

Table 3: Abliation study on SWaT.

Model	Prec.	Rec.	F1.
MGADN	99.50	76.95	0.87
MGADN w/o MGDA-UB(i.e MGADN-ALT)	98.45	69.44	0.84
MGADN w/o VAE head	98.97	72.56	0.81
MGADN w/o predicted head	77.89	58.99	0.67
MGADN w/o VAE head w/o Shared weighted layer(i.e GDN)	99.35	68.12	0.81
MGADN w/o predicted head w/o Shared weighted layer(i.e VAE)	74.75	57.79	0.65

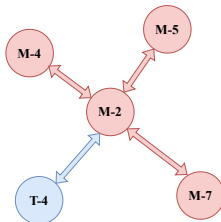


Figure 3: Sensor M-2’s neighbor learned by MGADN. Our model succeeded in capturing the relationship between M-7, M-4, M-5, which are highly correlated in reality. The potential relation between M-2 and T-4 is also revealed.

4.6. Other Interpretability

In this section, some visualization work on MSL will be completed to better interpret the reliability of MGADN.

When we chose sensor M-2 as an example, according to the topk-masked adjacent matrix, M-4, M-7, M-5 are regarded as the neighborhood of M-2. The description of the dataset tells us that there are actually connections between these sensors because they are sensors that receive the similar signals on Curiosity Rover. Potential connection between T-4 and M-2 are also revealed. Fig. 3 shows this situation.

Fig. 4 shows the variation of the unmasked adjacent matrix during the growth of epoch.

5. Conclusion

In this paper, we propose Multi-task Graph Anomaly Detection Network(MGADN), which is a two-headed time series anomaly detection approach. MGADN is able to capture relation between sensors automatically without any prior information. Multi-task optimization theory is also considered when tackling with losses of two different heads. Experiments on three real-world datasets proves its competence.

Future work may consider different kinds of reconstruction model or better graph embedding approach when learning graph structure. For instance, reconstruction models like

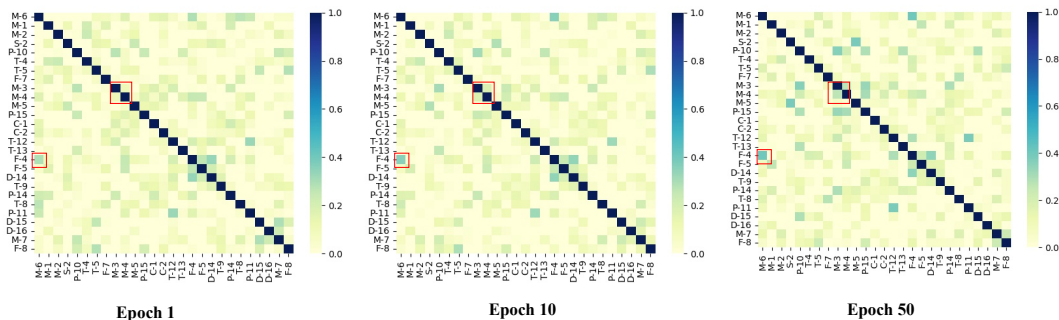


Figure 4: The adjacent matrix learned by MGADN without topk mask. As the number of epoch grows from 1 to 50, MGADN gradually capture the relationship between M-3 and M-4. Moreover, the relationship between sensors which are not intuitively connected is also learned such as F-4 and M-6.

GAN or LSTM-VAE may be introduced to multi-head framework. Instead of learning one-directed graph structure, we may also make it bi-directional via different graph embedding approach.

References

CC Aggarwal. Outlier analysis. data mining, 2015.

Chuahdhy Mujeeb Ahmed, Venkata Reddy Palleti, and Aditya P Mathur. Wadi: a water distribution testbed for research in the design of secure cyber physical systems. In *Proceedings of the 3rd International Workshop on Cyber-Physical Systems for Smart Water Networks*, pages 25–28, 2017.

Fabrizio Angiulli and Clara Pizzuti. Fast outlier detection in high dimensional spaces. In *European conference on principles of data mining and knowledge discovery*, pages 15–27. Springer, 2002.

George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. *Time series analysis: forecasting and control*. John Wiley & Sons, 2015.

Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29, 2016.

Ailin Deng and Bryan Hooi. Graph neural network-based anomaly detection in multivariate time series. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 4027–4035, 2021.

- Jean-Antoine Désidéri. Multiple-gradient descent algorithm (mgda) for multiobjective optimization. *Comptes Rendus Mathématique*, 350(5-6):313–318, 2012.
- Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric. *arXiv preprint arXiv:1903.02428*, 2019.
- Fazle Karim, Somshubra Majumdar, Houshang Darabi, and Shun Chen. Lstm fully convolutional networks for time series classification. *IEEE access*, 6:1662–1669, 2017.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Jitendra Kumar, Rimsha Goomer, and Ashutosh Kumar Singh. Long short term memory recurrent neural network (lstm-rnn) based workload forecasting model for cloud datacenters. *Procedia Computer Science*, 125:676–682, 2018.
- Dan Li, Dacheng Chen, Baihong Jin, Lei Shi, Jonathan Goh, and See-Kiong Ng. Mad-gan: Multivariate anomaly detection for time series data with generative adversarial networks. In *International Conference on Artificial Neural Networks*, pages 703–716. Springer, 2019.
- Wenkai Li, Wenbo Hu, Ning Chen, and Cheng Feng. Stacking vae with graph neural networks for effective and interpretable time series anomaly detection. *arXiv preprint arXiv:2105.08397*, 2021.
- Helmut Lütkepohl. Vector autoregressive models. In *Handbook of research methods and applications in empirical macroeconomics*. Edward Elgar Publishing, 2013.
- Aditya P Mathur and Nils Ole Tippenhauer. Swat: A water treatment testbed for research and training on ics security. In *2016 international workshop on cyber-physical systems for smart water networks (CySWater)*, pages 31–36. IEEE, 2016.
- Daehyung Park, Yuuna Hoshi, and Charles C Kemp. A multimodal anomaly detector for robot-assisted feeding using an lstm-based variational autoencoder. *IEEE Robotics and Automation Letters*, 3(3):1544–1551, 2018.
- Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*, 2017.
- Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018.
- Mei-Ling Shyu, Shu-Ching Chen, Kanoksri Sarinapakorn, and LiWu Chang. A novel anomaly detection scheme based on principal component classifier. Technical report, Miami Univ Coral Gables Fl Dept of Electrical and Computer Engineering, 2003.
- George Udny Yule. On a method of investigating periodicities in disturbed series, with special reference to wolfer’s sunspot numbers. *Philosophical Transactions of the Royal Society of London Series A*, 226:267–298, 1927.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.

- Gilbert Thomas Walker. On periodicity in series of related terms. *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character*, 131(818):518–532, 1931.
- Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang. Connecting the dots: Multivariate time series forecasting with graph neural networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 753–763, 2020.
- Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. *arXiv preprint arXiv:1709.04875*, 2017.
- Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*, 2014.
- Yu Zhang and Qiang Yang. An overview of multi-task learning. *National Science Review*, 5(1):30–43, 2018.
- Bin Zhou, Shenghua Liu, Bryan Hooi, Xueqi Cheng, and Jing Ye. Beatgan: Anomalous rhythm detection using adversarially generated time series. In *IJCAI*, pages 4433–4439, 2019.