

Gottfried Wilhelm Leibniz Universität Hannover
Fakultät für Elektrotechnik und Informatik
Institut für Verteilte Systeme

Master thesis
Informatics (M.Sc.)

Anomaly detection in streaming data using autoencoders

Student:	B.Sc. Bin Li
First Supervisor:	Prof. Dr. Eirini Ntoutsi
Second Supervisor:	Prof. Dr. Wolfgang Nejdl
Date:	July 18, 2018

Declaration of Authorship

I hereby certify that this thesis has been composed by me and is based on my own work, unless stated otherwise. No other person's work has been used without due acknowledgement in this thesis. All references and verbatim extracts have been quoted, and all sources of information have been specifically acknowledged.

Hanover, July 18, 2018

B.Sc. Bin Li

Contents

Abstract	1
1 Introduction	3
2 Related works	7
2.1 Batch model	7
2.1.1 Classical machine learning based approaches	7
2.1.2 Autoencoder-based batch approaches	8
2.2 Online anomaly detection over data stream	9
2.2.1 Classical online approaches	9
2.2.2 Autoencoder-based online approaches	10
3 Preliminaries	13
3.1 Definition of stream	13
3.2 Definition of anomalies	13
3.3 LSTMs	15
3.4 Autoencoders	16
4 Proposed model	19
4.1 Framework overview	19
4.2 LSTMs-Autoencoder	19
4.2.1 Encoder-decoder architecture	19
4.2.2 Online anomaly detection	21
4.3 Online learning	22
4.3.1 Updating dataset	22
4.3.2 Updating trigger	23
4.3.3 Model updating	24
5 Experimental setup	25
5.1 Datasets	25
5.2 Parameter tuning	26
5.3 Streaming data generator: Apache Kafka	27

6	Experimental results	29
6.1	Anomaly detection performance	29
6.2	Synthetic example	30
6.2.1	Reaction of concept drift	31
6.2.2	Model updating	31
6.3	Model retraining	32
6.3.1	Reaction of sudden and drastic concept drift	32
6.3.2	Reaction of serial concept drift	33
7	Conclusion	35

List of Figures

3.1	Data stream	13
3.2	Anomaly types	14
3.3	LSTM unit	15
3.4	Unfolded LSTM unit	16
3.5	Deep LSTMs	16
3.6	Autoencoder	17
3.7	LSTMs-Autoencoder	18
4.1	Online anomaly detection flowchart	20
4.2	Updating data buffer	23
5.1	Data stream Publisher-Consumer architecture	26
5.2	Kafka diagram	27
6.1	AUC comparison between stationary and adaptive models over stream	30
6.2	Run time statistics	31
6.3	Retraining effect on PowerDemand dataset	31
6.4	PowerDemand dataset online learning scores	32
6.5	PowerDemand dataset online parameter updating	32
6.6	SMTP+HTTP stream concept drift	33
6.7	ForestCover stream concept drift	34

List of Tables

3.1	Confusion matrix	14
5.1	Datasets information	25
5.2	Hyperparameters	27
6.1	Performance	29

Abstract

The high-volume and velocity data stream generated from devices and applications from different domains grows steadily and is valuable for big data research. One of the most important topic is anomaly detection for streaming data, which has attracted attention and investigation in plenty of areas, for example, the sensor data anomaly detection, predictive maintenance, event detection. Those efforts could potentially avoid large amount of financial costs in the manufacture. However, different from traditional anomaly detection tasks, anomaly detection in streaming data is especially difficult due to that data arrives along with the time with latent distribution changes, so that a single static model doesn't fit streaming data all the time. An anomaly could become normal during the data evolution, therefore it is necessary to maintain a dynamic system that adapt the changes. In this work, we propose a LSTMs-Autoencoder anomaly detection model for streaming data. This is a mini-batch based streaming processing approach. We experimented with streaming data that containing different kinds of anomalies as well as concept drifts, and the results suggest that our model can sufficiently detect anomaly from data stream and update model timely to fit the latest data property.

Chapter 1

Introduction

Anomaly detection attracts more and more attention in the data mining domain, and is widely used in the different applications, e.g. manufacture, e-commerce, internet etc. Successful anomaly detection in time avoids inconvenient and reduces maintenance expenditure in many scenarios like credit card fraud detecting, spam email recognition, machine condition monitoring, and could also be used as a preprocessing step to remove anomalies from datasets before machine learning tasks. There are already plenty of anomaly detection techniques proposed in previous literatures that solve this problem from variety perspectives, e.g. distance-based methods, clustering analysis, density-based methods etc.

There is no lack of anomaly detection approaches that perform good with respect to different kinds of data. Supervised approaches take anomaly detection as a binary classification problem of “normal” and “abnormal” instances, and all instance labels should be available in advance. The key difference to other classification problem is here the class labels are extremely biased to the normal class while anomaly only appears rarely. Instead of doing data augmentation on anomaly data for supervised approaches, unsupervised approaches are more direct solutions to this problem, which treat the instances that fit least to the majority as the anomalies. Furthermore, under many real-world situations, only partial labels are available, and therefore semi-supervised as well as one-class models are more efficient. Those models learn pattern from the partially labeled normal data, and scale anomaly likelihood according to the difference between an unseen pattern and the learned normal pattern.

However, most of the existing anomaly detection models are designed as batch models, which means, the entire training set should be available in advance. This becomes a shortcoming under today’s big data background. With the rapid development of hardware in the last decade, the situation of data acquisition and analysis has significantly been changed. Specifically, in the IoT applications, data are acquired from sensors attached to IoT devices, and arrive continuously and everlasting. In the beginning, no static entire training set is available for model initialization in the batch fashion. Besides, during data analysis, we should always consider the volume and velocity of data, which means,

on one hand, with traditional batch classifiers, the infinity data stream will lead to out of memory problem, on the other hand, streaming data usually comes with a high velocity that leaving the system few processing time. Optimally, the model should have only single look at each data point in the stream. In addition, the statistical property of data may also change over time, which is formally called ‘concept drift’. The model should always learn latest knowledge from the stream and update its identification of anomaly automatically, while anomalies could be temporally. After a data distribution change, an anomaly could possibly become normal in the new streaming context. Data distribution changes should not be classified as anomaly, in the meanwhile, though anomalies show up rarely over the stream, they should not be oversighted. To this end, an anomaly detection system for streaming data should be able to 1) be initialized with only a small subset, 2) process streaming data and make prediction in real-time, 3) adapt data evolution over the stream. 4) model should be able to deal with the biased data.

LSTMs are a kind of recurrent neural network that exhibit dynamic temporal behaviour for time series. As a neural network-based model, LSTMs are able to deal with high-dimensional and non-linear data directly. In the last decade, LSTM are used widely in time series prediction and text prediction. LSTMs-based autoencoders are also applied to sequence to sequence problems, e.g. language translation, time series data embedding. Deep LSTMs have been shown good performance in capturing hierarchical information of time series like separation of sentences. Recently, LSTMs-based autoencoders are also used for time series anomaly detection in order to capture the temporal information between data points. For example, Malhotra et al. introduced an autoencoder based anomaly detection approaches in [1],[2], and achieved good performance in multiple time series dataset. However, in those researches, they still assume that the whole datasets are available beforehand and work on static data. Also, the aforementioned on-line learning difficulties are not taken into consideration. Hence, we enhanced this kind of LSTMs-Autoencoder based static anomaly detection approaches with the online learning ability by appending incremental model updating strategies.

Neural networks, including autoencoders, are normally used in batch fashion, namely the whole training set is available, and trained with backpropagation. Under online setting, only small subset accumulated data from stream are available for model initial training, which may lead to getting a suboptimal model. The precondition of the online LSTMs-autoencoder is that the accumulated initialization set is enough to train a convergent model. The difficulty is to detect when model should be updated and what data should be used for model updating. The short-term changes of data distribution should not cause model variation, while permanent concept drifts should trigger model updating as soon as possible.

In this paper, we introduce a novel and robust incremental LSTMs-Autoencoder anomaly detection model, which designed specifically for time series data in a streaming fashion

using Long Short-Term memory (LSTM) units, with also online learning ability for model updating. For each accumulated mini-batch of streaming data, the autoencoder reconstructs it with previous knowledge learned from normal data. Anomalies (never used for training) are supposed to cause significant larger reconstruction error than normal data. In addition, the model is able to update itself when data distribution changes are detected. The LSTMs-Autoencoder is experimented with different data streams, and the experiment results show its ability in detecting anomaly from streaming data and is able to adjust itself with different kinds of concept drifts by model online updating.

The rest of this paper is organized as follow. In chapter 2, we collected previous works on anomaly detection and their shortcomings. We also refer previous works on incremental neural network. In chapter 3, define the problem formally. In chapter 4, we propose our method for streaming data anomaly detection and discuss the advantage over previous works. In chapter 5, we describe the datasets used for experiments and the experimental set-up. In chapter 6, we present our experimental results. And in chapter 7, we summarize the work and discuss of success and deficiency.

Chapter 2

Related works

In this section, we present a survey on previous works in anomaly detection. Both batch models and online models are listed with respectively classical machine learning approaches and neural network as well as autoencoder based approaches. And for the online case, we also survey works with neural network updating approaches.

2.1 Batch model

2.1.1 Classical machine learning based approaches

As an important component of data mining and machine learning, anomaly detection has been investigated using plenty of efficient models. When talking about anomaly detection, the most intuitive solutions are detection of outliers from a dense cluster, or to find data points that have obvious different property as their neighbors. Considering the lack of label and extremely imbalanced dataset, unsupervised approaches are more widely used in practice, for example Local Outlier Factor (LOF).

In anomaly detection, the LOF is a common density-based approach. LOF shares some concepts with DBSCAN such as ‘core distance’ and ‘reachability distance’, in order to estimate local density. Here, points with substantially lower local density than their neighbors are considered as anomalies. LOF shows competitive performance in many anomaly detection tasks, especially when dealing with data with unevenly density distribution. However, when getting a numerical factor from LOF model as result, it is actually hard to define a threshold automatically for the judgement of anomaly.

Because the anomalies appear rarely in the dataset, and occur usually in novel ways, it is expensive to label and hard to learn all kinds of anomalies during the training phase, so unsupervised models are commonly used. There are also a batch of semi-supervised or one-class anomaly detection models. The intuitive difference between anomaly detection and binary classification problem is the obvious fewer positive class data (anomalies). A typical one-class model is the One-class Support Vector Machine (OCSVM).

As an semi-supervised one-class classifier, OCSVM takes only normal data as input, and generates a decision surface to separate them from the anomaly data. By analyzing anomalies, the datasets are always bias to the normal part, and anomaly appear only rarely. So, this kind of one-class classifiers avoid making balance between the two classes. Besides, they also take advantage of classical support vector machine by using kernel methods, they can also deal with linearly not separable data. However, in the meantime, choosing a proper kernel becomes a hard problem for OCSVM. A suboptimal kernel function can seriously impact the performance.

Although classical machine learning approaches can handle most of the normal anomaly detection, there is still a lack of pervasive models that fit different kinds of data characters. Moreover, only few of those approaches could be directly or after some modification used for time series and streaming data, while they ignore the temporal dependency between samples.

2.1.2 Autoencoder-based batch approaches

LSTMs-Autoencoders are originally widely used for text generation because the LSTMs are able to capture the contextual dependency between words and sentences. Text data are usually embedded into vector as input data of autoencoder. And the tasks are either to generate temporal relevant text on the decoder side or to learn text representation from the hidden layer [1]. Similar problem appears also in time series mining due to the temporal dependency of values between timestamps. So, in later works, LSTMs-Autoencoder are also used for time series data.

Sutskever et al. [13] proposed a deep LSTMs-based sequence-to-sequence model for language translation. In their work, the deep LSTMs encoder takes single sentence as input, and learn a hidden vector with fixed length, then a different LSTMs decoder decodes the hidden vector to the target sentence. As a translation task, they found that this encoder-decoder architecture can capture long sentences and sensible phrases, especially they achieved better performance with deep LSTMs in compare with shallow LSTMs. In addition, a valuable found is, reversing the order of words in the input sentence makes the optimization problem much easier and achieved better performance. The LSTMs based model outperforms non-LSTMs model on the long input sentence cases (more than 35 words) since its long-term memory ability.

Li et al. [7] did similar research on long paragraph text and even entire document generation using LSTMs-autoencoders. Their main contribution is the hierarchical sentences representation. The model learns word-level, sentence-level and paragraph-level information with each respectively a LSTMs layer, so that the model captures very long-term temporal information. Moreover, they introduced an attention based hierarchical sequence-to-sequence model that connect the most relevant parts between encoder and decoder, for example, the words around a final punctuation. They experiment with docu-

ments over 100 words, the results show that hierarchical and attention-based hierarchical LSTMs learns even better long-term temporal information than standard LSTMs-based encoder-decoder models.

As autoencoders achieve great successes in text data and speech processing, they are also used for time series anomaly detection. These models train autoencoders with only normal data while take anomalies as unknown patterns. Then the autoencoders are only able to reconstruct normal patterns, then large reconstruction error indicates anomaly. An early work [11] used vanilla autoencoder to detect abnormal status of the electric power system. In order to capture temporal information, they applied sliding window on the raw data as input. As anomaly scoring method, they evaluated each sliding window with respect to their reconstruction error. As some measures in the autoencoder output vectors that are more sensible to anomalies than others, they use the average absolute deviation of reconstruction error as anomaly score. And the anomaly threshold is chosen by large amount of experiments over normal data.

Another important reason of using autoencoder for anomaly detection is its ability of dealing with high-dimensional data. Sakurada et al. [12] experimented with time series data that consist of 10-100 variables without linear correlation. Comparing with reconstruction using PCA or Kernel PCA techniques, the autoencoder reconstruction error can easily recognize anomalies.

In further researches, Malhotra et al. [10][8] developed the application of LSTMs-autoencoder from sequence learning into anomaly detection problem. They proposed a stacked LSTMs model to learn high level temporal patterns. They show that LSTMs outperform normal RNNs-based anomaly detection model and avoid the gradient vanishing problem thanks to the gate architecture inside the LSTM unit. They also detect anomaly based on the reconstruction error. The scoring function is based on the parameters of an estimated normal distribution of a validation set. Their experiments show that the model performs good in variety kinds of datasets. A variation of this model [9] has been proved that achieves better performance in the anomaly detection tasks, while they tell that using a constant as input of decoder instead of read time series value improves the performance of model.

2.2 Online anomaly detection over data stream

2.2.1 Classical online approaches

Recently, streaming data mining attracts more and more attention, and many efficient classical models are modified to fit the online learning property. Cui et al. [2] proposed online anomaly detection approach with grid-based summarization and clustering algorithm, which is able to detect anomaly immediatly when data arrives. For the passed streaming data, they used summarization algorithm to reduce the run time and memory

consumption over stream. And they give anomaly scores to instances to describe concrete anomaly degree. The model shows good performance in KDD dataset for network attack detection. However, they consider only isolated data point without any temporal information measuring. Another streaming data anomaly detection framework by Tang et al. [15] used sliding window to involve the contextual dependency and temporal changes in the stream. The anomaly detection algorithm is a modified version of LOF while LOF's high time complexity can not be directly employed to streaming data. They use comentropy to filter out most normal windows, and feed the rest to LOF.

2.2.2 Autoencoder-based online approaches

Zhou et al. [16] proposed an online updating approach for denoising autoencoders by modifying the hidden layer neurons in order to deal with the non-stationary streaming data properties. The basic idea consists two steps, adding hidden layer neurons to capture new knowledge, and merging hidden layer neurons if information is redundant. Their experimental result shows comparable or better reconstruction result than non-incremental approaches with only few data used during initialization. And they show that their incremental feature learning methods performs more adaptively and robustly to highly non-stationary input data distribution.

Dong et al. [3] proposed a 2-step anomaly detection mechanism with incremental autoencoders. They implemented the model with ensembled autoencoders in multithreads in order to leverage parallel computing when large volumes of data arrive. In the 2-step mechanism, they check anomaly in the first step and verify anomaly data with previous and subsequent data (to distinguish between anomalous state and concept drift) in the second step to reduce false-positive rate in anomaly detection. In the experimental results, they show that their model outperforms commonly used tree-based streaming anomaly detection models especially when concept drift presents. And they speed up the online processing with mini-batch learning and online learning in multithreads.

Ghazikhani et al. [4] introduced an online neural network model for streaming data towards to the two major problems of online learning, concept drift and imbalanced classes. In term of concept drift, they applied a forgetting function that weights recent instances to navigate the model to the drifted model, so that the model always learns pattern from latest data. Besides, for class imbalance, they proposed an error function for two-class imbalance problem with the basic idea that the error function generating higher error signals for instances in the minority class.

Kochurov et al. [6] designed incremental learning framework for deep neural networks based on Bayesian inference. They argued that, naïve deep learning approaches for incremental learning applies Stochastic Gradient Descent (SGD), which intent to keep previous learned model remembered, and enhanced with current batch of new data. However, by SGD, the neural network model is likely to converge to the local optimal

of the latest batch of data without preserving the previous knowledge. Their Bayesian framework estimates the posterior distribution over the weights of the model in the condition of previous knowledge and use the Bayesian rule to sequentially update the posterior distribution in the incremental learning.

In this work, we implement a LSTMs-autoencoder based incremental streaming data anomaly detection model. The LSTMs-autoencoder is close to the model introduced by Malhotra et al. [8], and we design an online model updating strategy as well as the updating data buffers used for model updating, which collect fresh knowledge from the last seen instances.

Chapter 3

Preliminaries

3.1 Definition of stream

Assume that a set of devices or data warehouses generate data continuously (here we only take about numerical data). A data stream DS is defined as

$$DS = \{(x_1, t_1), (x_2, t_2), \dots, (x_T, t_T), \dots\} \quad (3.1)$$

where x_T is the instance arrived at timestamp t_T and represented by a multi-dimensional vector.

In order to feed the streaming data to LSTMs, the stream is accumulated to windows and batches as Figure 3.1. Every T instances are accumulated to a window as a single LSTM input. And MB is the predefined batch size, each batch contains MB windows.

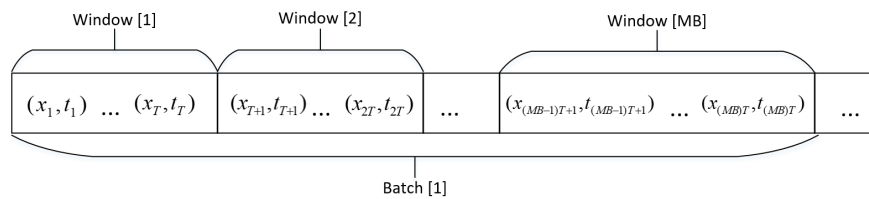
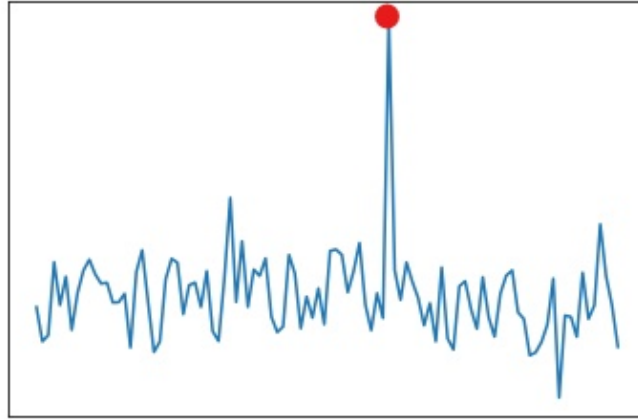


Figure 3.1: Data stream

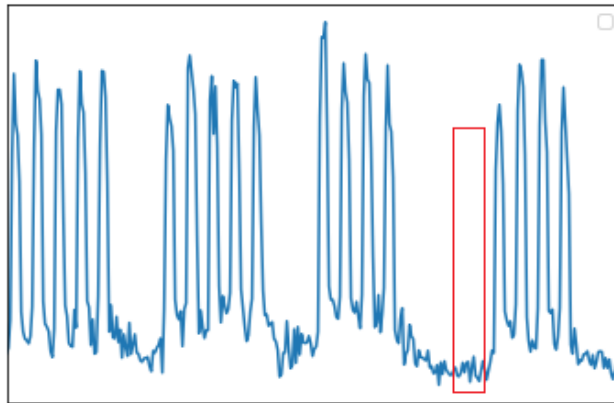
3.2 Definition of anomalies

Pointwise: A data point (instance) is anomalous (e.g. Figure 3.2a on the next page) if this point is distant from other observations according to some specific measurement metrics. This is used in fine-grained anomaly detection tasks, that need to find out every single anomalous instance, e.g. credit card fraud detection, spam email detection.

Window-based: Sometimes, a data point is apparently normal, but this point, or potentially together with its neighbors violates the overall periodicity or other character of the time series, we also treat them as anomaly, which is called window-based anomaly or contextual anomaly, e.g. Figure 3.2b.



(a) Pointwise anomaly



(b) Window-based anomaly

Figure 3.2: Anomaly types

In the anomaly detection experiments, abnormal data is the object of study, therefore we take the anomaly as positive class and normal data as negative class.

		Actual value	
		Normal	Abnormal
Prediction	Normal	TN	FN
	Abnormal	FP	TP

Table 3.1: Confusion matrix

The target is to achieve higher true positive rate (equation 3.2, true alarms) while remain lower false positive rate (equation 3.3, wrong alarms). The evaluation metric is Area Un-

der the Curve (AUC), where the curve is receiver operating characteristic curve, and is created by plotting the true positive rate against the false positive rate at various threshold settings. The range of AUC is between 0 and 1, 1 is the optimal result.

$$TPR = \frac{TP}{TP + FN} \quad (3.2)$$

$$FPR = \frac{FP}{FP + TN} \quad (3.3)$$

3.3 LSTMs

Recurrent neural networks(RNNs) are widely used for speech, video recognition and prediction due to its recurrent property that captures the temporal dependency between data points in compare with other feed forward networks. However, the volume of RNN's memory is limited, and vanishing gradient is also a difficulty by training RNNs. Therefore, the long short-term memory networks (LSTMs) are a kind of reinforced RNN that are able to remember valuable information in arbitrary time interval. A LSTM network is a recurrent neural network with neurons being LSTM units. Figure 3.3 shows a classical structure of a LSTM unit. LSTMs are able to capture long-term memory while there are a forget gate and a update gate in the LSTM unit, that select necessary previous information and new coming information according to the input data at each time step. The information is transferred to the next step along with the cell state. Besides, each LSTM units also output its value respectively.

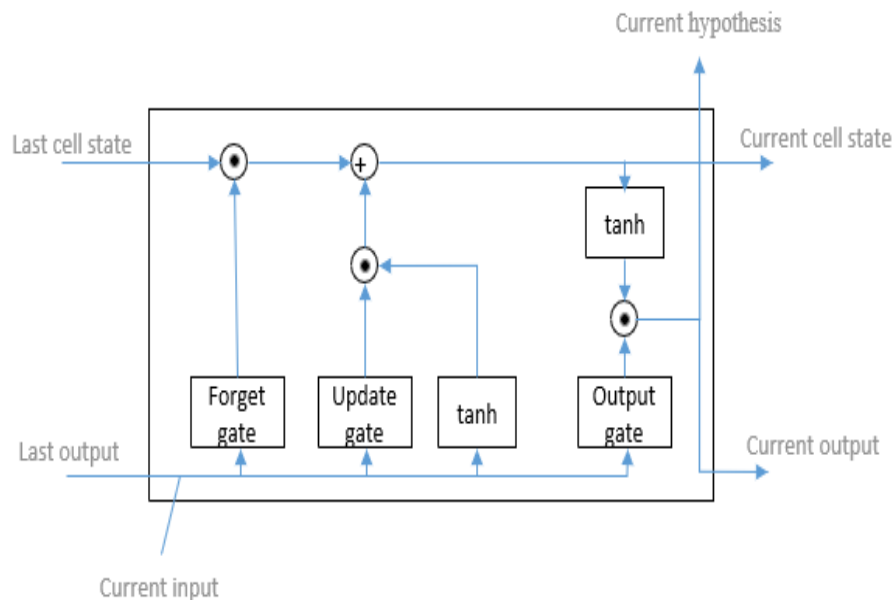


Figure 3.3: The LSTM unit

A LSTM unit can be unfolded over time, as shown in Figure 3.4 on the following page. The LSTM unit takes a data window as input (one instance at each time step). Therefore,

the LSTM unit extracts useful and drop useless temporal information from the window.

Deep LSTM RNNs are built by stacking multiple LSTM layers. Note that LSTM RNNs are already deep architectures in the sense that they can be considered as a feed-forward neural network unrolled in time where each layer shares the same model parameters. It has been argued that deep layers in RNNs allow the network to learn at different time scales over the input[5]. Figure 3.5 is a example of stacked deep LSTM neural network, there are 3 LSTM layers, each can be unfolded into 5 time steps, so the LSTMs take a window in length 5 as input and the output is in same size.

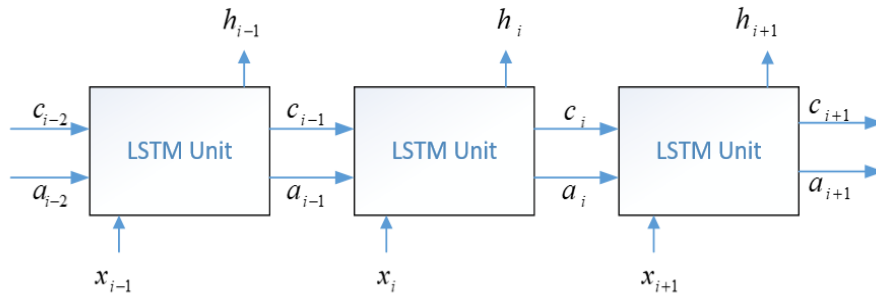


Figure 3.4: Unfolded LSTM unit

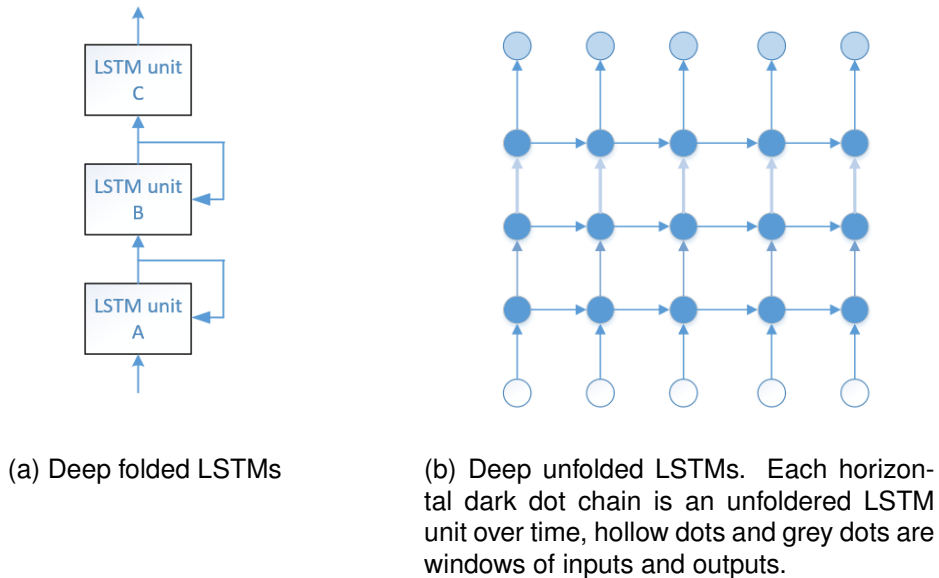


Figure 3.5: Deep LSTMs

3.4 Autoencoders

An autoencoder (Figure 3.6 on the facing page) is an artificial neural network with symmetrical structure. Normally an autoencoder has at least one hidden layer that consists

of less neurons than input and output layers. And the basic aim of autoencoders is to reconstruct its own input and learn a lower dimensional representation (encoding) of input data in the hidden layer. Moreover, the autoencoders are also used for anomaly detection by measuring the reconstruction error between inputs and predictions. Normally the component between input layer and hidden layer is called encoder, and the symmetrical component between hidden layer and output layer is called decoder. For input χ , the objective function is to find weight vectors for encoder and decoder to minimize the reconstruction error (Equation 3.4).

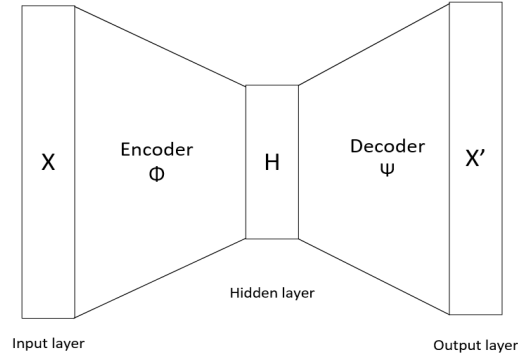


Figure 3.6: Autoencoder

$$\begin{aligned}
 \Phi : \chi &\rightarrow H \\
 \Psi : H &\rightarrow \chi \\
 \Phi, \Psi &= \underset{\Phi, \Psi}{\operatorname{argmin}} \|\chi - (\Psi \circ \Phi)\chi\|^2
 \end{aligned} \tag{3.4}$$

LSTMs-autoencoder has the same encoder-decoder architecture, while the neurons are LSTM units and connected in the way described in section 3.3. Figure 3.7 on the next page is a basic LSTMs-based autoencoder architecture with single LSTM layer on both encoder and decoder side. Our incremental LSTMs-autoencoder is based on this structure. The model takes window with length T as input (one instance at each step). The cell state carries sequence information and is passed through LSTM unit over time. When the encoder reaches the last encoder state, namely ET in Figure 3.7b on the following page, its cell state is actually the fix length embedding of the input window, and will be copied to the decoder as initial cell state of decoder, so that the input information is also transferred to the decoder. And the decoder predict the window in reversed order in order to make the optimization problem easier. To be notice is, different from aforementioned deep LSTMs in section 3.3, the encoder outputs at each time step are not directly used as inputs of decoder, while between the encoder and decoder is actually not the same logical connection as stacked LSTMs. Here, the outputs of encoder are ignored, and there are different works contributes to the research of decoder inputs. Cho et al. [1] feeds the input sequence to the decoder for a learning phrase representation task, Malhotra et al. [8] feed to decoder LSTM unit at each time step the value of previous time step as input, and in a extended work [9] they feed the decoder always a constant vector

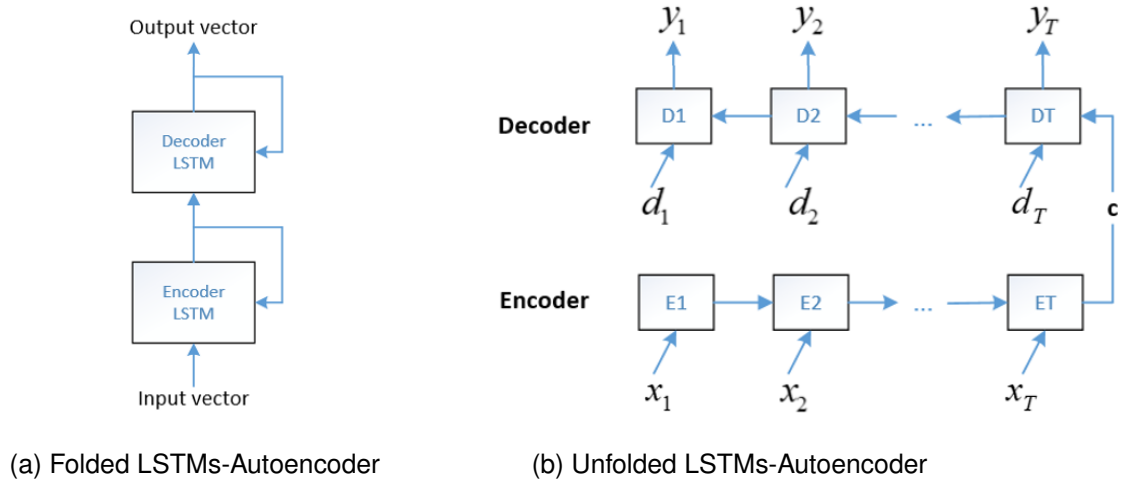


Figure 3.7: LSTMs-Autoencoder

for an anomaly detection task, because the final cell state already carries all relevant information to represent the input window. In our model, we feed the decoder a constant vector.

Chapter 4

Proposed model

4.1 Framework overview

The proposed model is a full flow from data stream generation, anomaly detection with autoencoder-based model and online model incremental updating. Figure 4.1 on the next page shows the general pipeline. The first received batches of streaming data are used for decision of model hyperparameters and the model initialization. Hyperparameters includes the hidden layer size, batch size, input window length as well as the number of epochs. Once the hyperparameters are determined, an autoencoder will be constructed and initialized with random weights. A subset of the streaming data is used for model initialization (only normal data used for training). Furthermore, the model is used for online anomaly detection, and evaluated based on the labels provided by experts. After evaluation, the data windows with bad performance are collected as hard examples in buffers for updating. Model will be updated when the updating condition is triggered. As shown in Algorithm 1, if a batch of streaming data is available, the model will start do prediction, evaluation, and check whether current window is useful to store for later updating. If so, the window of data will be appended to the updating buffer, with the instance order not been destroyed. As a consequence of anomalies' rare appearance, we keep all seen anomalous windows for determination of anomaly score threshold during updating.

4.2 LSTMs-Autoencoder

4.2.1 Encoder-decoder architecture

The LSTMs-Autoencoder consists of two LSTM units, one as encoder and the other as decoder. The encoder inputs are fix length vectors with shape $\langle MB, T, D \rangle$, where MB is the number of data windows contained in a mini-batch, T is the numbers of data points within each data window, and D represents the number of data dimensionality. Here, MB and T are learned as hyperparameter in the initialization phase. And on the decoder side, it is supposed to output exactly the same format data vector for each mini-batch. The LSTM unit copies its cell state for itself as one of the cell input at next timestamp. At the last timestamp of encoder, the cell state of LSTM unit is the hidden representation

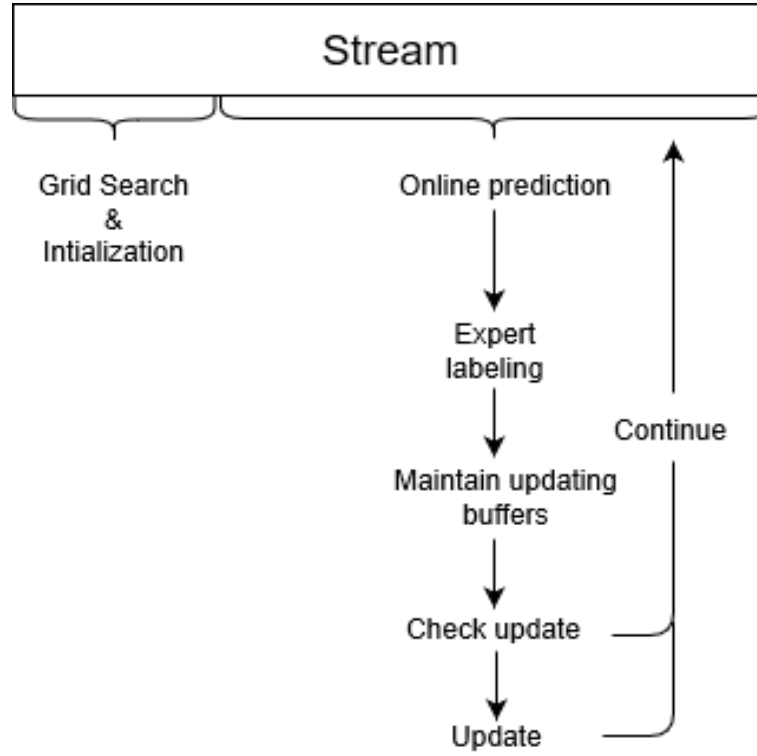


Figure 4.1: Online anomaly detection flowchart

of the input data vector and copied to the decoder unit as initial cell state, so the hidden information can be passed to the decoder. The size of hidden layer representation vector, namely the size of cell state is another hyperparameter need to be learn in the initialization phase. The larger the hidden vector, the more information can be captured during the process, so it is a feature highly depends on the data. Similar to previous study [13], we also train the encoder and decoder with time series in reverse order. For example, if the input data fragment are data points from timestamp t_1 to t_2 , then the decoder will predict data point at t_2 at first, and then back to t_1 step by step, while this trick makes the gradient escarpment between last state of encoder and first state of decoder smaller and easier to learn.

In order to let the whole process happen online, the model initialization also utilizes streaming data. Once a small subset of streaming data is available, hyperparameters are learned, and then another dataset that consists only of normal data is collected from stream used for training. Assume that once an anomaly detection task is determined, the anomalous state is explicit defined and a subset of anomalous data is available for model initialization. We split the normal data into four subsets, N_1 for hyperparameters tuning, N_2 for model training, N_3 for early stopping, and scoring parameters learning, N_4 for testing. And abnormal data are split into two subsets, A_1 for decision of anomaly score threshold, A_2 for testing.

Algorithm 1: OnlineAnomalyDetection

```

input: normal buffer size: SN, performance threshold: P

needUpdating = False;
normalBuffer = [ ];
abnormalBuffer = [ ];

while True do
  if len(normalBuffer) >= SN and len(abnormalBuffer) != 0 then
    update(normalBuffer[-SN:], abnormalBuffer);
    normalBuffer = [ ];
    abnormalBuffer = [ ];
  else
    data, labels = getBatchData();
    pred = predict(data);
    result = evaluate(pred, labels);
    foreach window in data do
      if 'anomaly' in label[window.index] then
        abnormalBuffer.append(window)
      end
      if AUC(result[window.index]) >= P then
        continue;
      else
        if label == 'normal' then
          normalBuffer.append(window);
        else
          continue;
        end
      end
    end
  end
end

```

4.2.2 Online anomaly detection

The autoencoder reconstructs the input with its knowledge of normal data, so if the input data contains anomalies, the reconstruction error will be obviously larger due to the lack of anomalous knowledge. For input $X^{(i)}$, the reconstruction error is

$$e^{(i)} = \left| X^{(i)} - X'^{(i)} \right| \quad (4.1)$$

similar to [8], the reconstruction error of data points in N_3 is used to estimate the parameters μ and Σ of a normal distribution $\mathcal{N}(\mu, \Sigma)$ using maximum likelihood estimation. The anomaly score for a point $x_t^{(i)}$ is defined as

$$a^{(i)} = (e^{(i)} - \mu)^T \Sigma^{-1} (e^{(i)} - \mu) \quad (4.2)$$

During the initialization phase, an anomaly score threshold τ is also learned using N_3 and A_1 as

$$\tau = \operatorname{argmax} AUC(a(N_3), a(A_1)) \quad (4.3)$$

The anomaly score of every instance in a window is compared with τ , and values over τ are predicted as anomalies. If a window contains more than τ_N anomalous values, this window is predicted as anomaly. τ_N is determined by dataset study.

4.3 Online learning

However, if the model is utilized for streaming data, the autoencoder will possibly become outdated because of the relative small and simple initialization dataset and concept drift that happened along with the time. So the update of model is indispensable. In this section, we introduce the updating strategy of the LSTMs-Autoencoder.

4.3.1 Updating dataset

Once the LSTMs-Autoencoder is initialized, it is ready for online prediction. There is a multi-thread setting in the online learning architecture. A sub-thread collects data instances continuously from the stream, and in the meantime, the main-thread works on real-time anomaly detection as long as mini-batches of data is provided by the sub-thread. For each single window in the mini-batch, every instance is reconstructed and the anomaly score is calculated using Equation (4.2) on the preceding page. The system maintains two data buffers for updating (Figure 4.2 on the next page), one for normal data windows, and the other one for anomalous windows. Considering the fact that a well mastered window leads to lower reconstruction error, and higher error indicates new features in the data, we can measure this reconstruction error level by the predefined normal distribution on reconstruction error. After each batch, the label for each data window is determined by experts. We predefine a performance threshold PN for normal data. Normal data windows that containing more than PN over τ_N instances are regarded as not good mastered and will be appended into the normal buffer for updating. As anomalies appear rarely in the stream, we collect all anomalous windows in the abnormal buffer for score threshold determination during updating.

Because the out-of-date buffer not might be collect from previous concept drift time period, and not benefits to current updating, we maintain the retain buffers with a queue structure, so that only a specific amount of most fresh data can stay in the buffer. To this end, once a updating process is triggered, only not well mastered fresh normal data are used for updating.

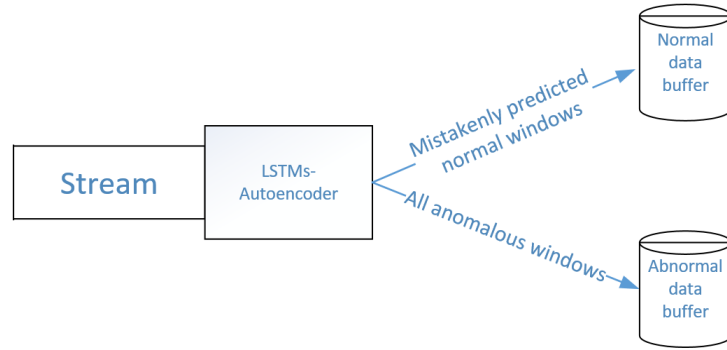


Figure 4.2: Updating data buffer

4.3.2 Updating trigger

During the online processing, if the system detects that the model doesn't fit the current data any more, then the model updating is triggered and done with the latest collected data in the two buffers. During experiments we found that, anomalies only appears rarely in the stream, so it often happens that the model need updating to fit the latest data, but still lack of anomaly data in the buffer to update the anomaly score threshold. To this end, we trigger the updating so long as the anomaly buffer is not empty. If the anomaly data are not enough to make up a single batch, then we duplicate the windows in anomaly buffer until buffer fits the batch size. In case of the normal buffer reaches a predefined size and anomaly buffer is not empty, the model is updated in a sub thread while the main thread keeps processing the stream.

The updating trigger strategy depends on the buffer size. This approach is suitable for larger, relative stationary data, while even concept drift happens, large amount of data arrives quickly to enrich the updating buffers, and trigger updating in time. And this approach highly depends on the performance threshold that decides when a data window from stream should be appended to the buffers, namely, updating is not directly depends on the real-time prediction performance.

Another updating trigger strategy is designed for smaller data set, where the waiting time of updating buffer full might be long after concept drift happening. During the waiting time, there can be other concept drifts, and the prediction performance is suboptimal during this time period. So, for smaller data sets, the updating trigger should directly relate to real-time performance. A simple way is, compare the batch performance with the first batch after last updating or streaming beginning. The reason is, the model is only updated with normal data, therefore every updating brings new knowledge to the model, and improve the performance, so the batch performance should at least same as or better than the first batch performance, otherwise it indicates concept drift.

4.3.3 Model updating

Once the updating process is triggered, the model will be updated using data from the buffers. Windows of normal buffer are divided into updating training set and updating validation set. Once the online phase starts, the LSTMs-Autoencoder is loaded into memory, and further model updating are all done in memory. The updating is a continuation of the initialization or previous updating (last check point) with identical data format. Parameters μ , σ as well as anomaly score threshold are learned from the updating validation set and anomaly buffer data. The parameters μ and σ are the mean and variance (or covariance for multivariate data) of reconstruction error estimated by normal validation set during training. So we learn new parameters in the updating using normal validation set as well.

Algorithm 2: update

input: nBuf, aBuf

trainSet, valSetN = split(nBuf);

valSetA = aBuf;

train(trainSet);

$\mu, \sigma, \text{threshold} = \text{getParameters}(\text{valSetN}, \text{valSetA});$

Chapter 5

Experimental setup

5.1 Datasets

We use 5 datasets in our experiments, PowerDemand, SMTP, HTTP, SMTP+HTTP and ForestCover. Those are widely used datasets in the streaming data mining area [8][3][14]. Statistical features are listed in Table 5.1. PowerDemand is a small univariate time series that records the power demand over a period of one year. Weekdays' power demand is higher than weekends' and daytime is higher than nights, power demand of special days (e.g. festivals) are abnormal. The aim is to find out such contextual anomalous weeks instead of single strange data point. The experimental data is an subsampling of original set. We demonstrate a synthetic example with visualization using this dataset while the trends and anomalous states are relative obviously. SMTP, HTTP, SMTP+HTTP are streaming anomaly data extracted from KDD Cup 99 dataset. According to Tan et al. [14], HTTP contains sudden surges of anomalies and SMTP does not, but possibly exhibits some distribution changes within the stream. Because of the difficulty to point out where the distribution changes occur in the stream, the HTTP+SMPT dataset is derived by connecting SMTP and HTTP, so that a distribution change is occurred when the communication protocol is switched. The ForestCover dataset is from the UCI repository, which contains 7 kinds of forest cover types. Similar as Dong et al. [3], we defined the smallest class Cottonwood/Willow with 2747 instances as anomaly, and the rest 6 classes as normal class with distribution changes.

Table 5.1: Datasets information

Dataset	Size	Dimensionality	Anomaly proportion(%)
PowerDemand	35 040	1	2.19
SMTP	96 554	34	1.23
HTTP	623 091	34	0.65
SMTP+HTTP	719 645	34	0.73
ForestCover	581 012	7	0.47

We separate each dataset into initialization set and streaming set, both contain normal and abnormal data. Further, the initialization set is divided into

- GS(n&a): for grid search
- INIT(n): for model initial training
- PL(n&a): for model parameter learning
- TEST(n&a): for testing during initialization

where “n” represents normal data and “a” represents abnormal data. And the streaming set is published to Kafka to generate data stream (Figure 5.1).

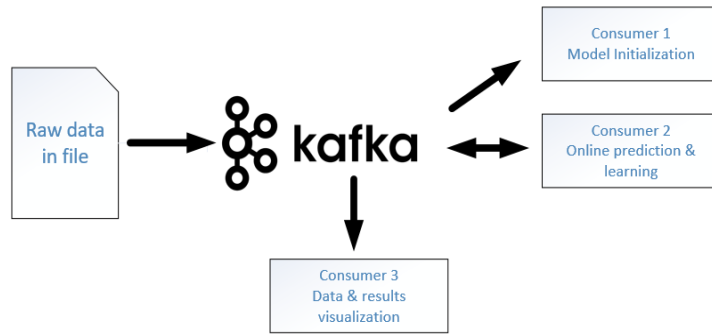


Figure 5.1: Data stream Publisher-Consumer architecture

For each dataset, we use around 20% data for initialization and the other half for online prediction. Each subset used for training and prediction are preprocessed locally in batch, in order to scale them into $[0, 1]$ to fit the LSTM activation function.

5.2 Parameter tuning

For each dataset, we carry out a grid search step to tune the model hyperparameters. Here we try multiple combinations of window length and hidden size for each data set. The grid search set GS contains 5% -15% anomalies. Because of the uncertainty of the random neural network weight initialization, we do each experiment 10 times and take the average result to reduce the impact. To be noted that during every divisions, the consistency of streaming data is persisted, or in other words, there is no random sampling during the process. A good model should make the reconstruction error as large as possible in order to make the classification easier. Given dataset D , ‘win’ is the input window, the average reconstruction error of D is given by Equation (5.1). The target function of grid search is given by the ration of average reconstruction error of abnormal and normal test grid search set GS (Equation (5.1)).

$$ARE(D) = \frac{1}{T} \sum_{t=1, win \in D}^T (input_{t,win} - output_{t,win}) \quad (5.1)$$

$$REratio = \frac{ARE(G_a)}{ARE(G_n)} \quad (5.2)$$

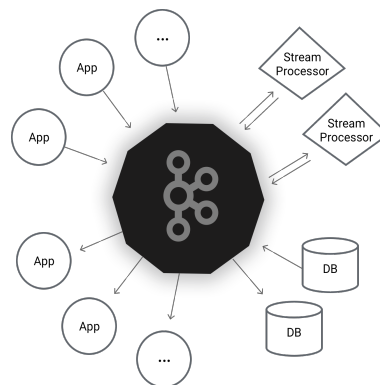
As a result of grid search, the hyperparameters for different datasets are listed in Table 5.2.

Table 5.2: Hyperparameters

Dataset	Window length	Hidden size	#Grid Search instance
PowerDemand	80	15	1 000
SMTP	10	15	5 000
SMTP+HTTP	10	15	5 000
HTTP	30	35	10 000
ForestCover	10	25	10 000

5.3 Streaming data generator: Apache Kafka

We utilize Apache Kafka as the streaming platform. Kafka is a widely used Publish/-Subscribe architecture streaming system. It is different from classical message queue technique with its fault tolerant, durable and large capacity properties. Different application or database can publish data to a specific topic of Kafka (topic is the data category mechanisms used in Kafka), and other processors can consume data from this topic (Figure 5.2). In our experimental setting, the data source is static databases, Kafka generate real-time data stream pipeline as data source publish records to the experiment topic, and furthermore the stream of records will be consumed by different consumers like our analysis model, visualization model etc. This configuration can be easily scaled up to more complicated and demanding real world use cases. Each record in the Kafka stream pipeline is in the form of [Key, Value, Timestamp], where keys are used for positioning and values carry the data record.

Figure 5.2: Kafka diagram¹

¹From: <https://kafka.apache.org/>

Chapter 6

Experimental results

6.1 Anomaly detection performance

With parameters learned from Section 5.2 on page 26, autoencoders are trained for each dataset with the beginning of streaming data. The anomaly detection performance is indicated by AUC. For each dataset, we compare the AUC of online phase that without and with continuously model and parameter updating (Table 6.1). For the PowerDemand dataset, retrain trigger depends on the batch performance. And for the rest datasets, retraining only triggered when retrain buffers are full. The retraining brings overall performance improvement on all datasets comparing to stationary models. Especially in the SMTP+HTTP dataset, the stationary without learning concept drifted knowledge performs clearly worse than the model with updating.

Table 6.1: Performance

Dataset	AUC(without retraining)	AUC(with retraining)	#retrain
PowerDemand	0.91	0.97	2
SMTP	0.94	0.98	2
HTTP	0.76	0.86	2
SMTP+HTTP	0.64	0.85	3
ForestCover	0.74	0.82	8

In order to compare the performance with and without retraining, after each model updating, we calculate the AUC value for each specified time period. As Shown in Figure 6.1 on the following page, the x-axis represents the periods, for example, before first retraining (shown as P1 in each subplot), between first and second retraining (P2), and so on. For each dataset, we compare the AUC value of stationary model (trained with only initialization set) and adaptive model (online updated). For most cases, the adaptive models outperform stationary models, which shows the models profits from the knowledge updating over streaming data. To be notice that the SMTP+HTTP set contains sudden concept drift around P3, which leads to a sharp decline of the stationary model. In the meantime, the adaptive model is only slightly influenced by the mixed knowledge at P3 but keeps outstanding performance when the stream switches to HTTP side.

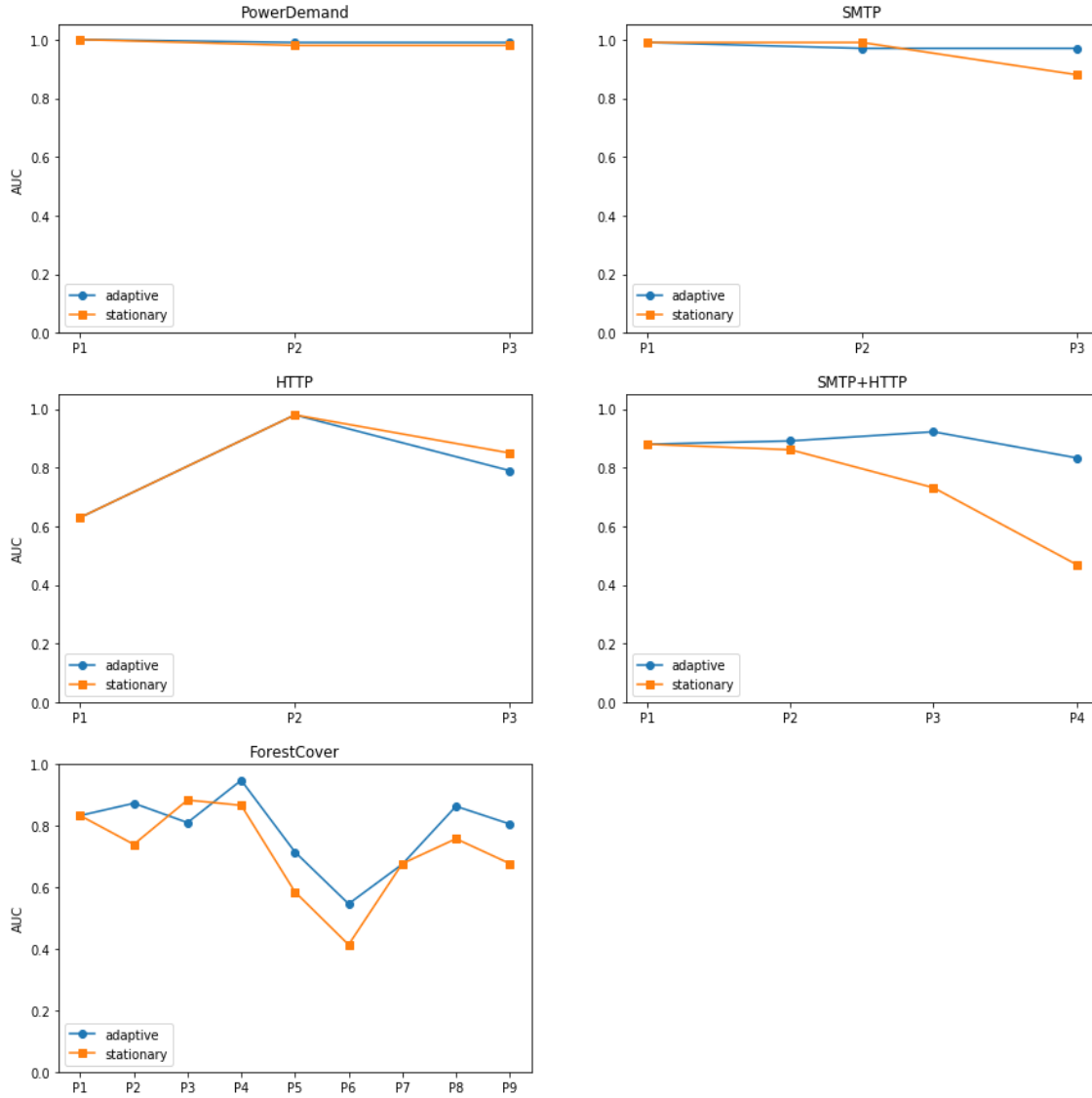


Figure 6.1: AUC comparison between stationary and adaptive models over stream. The x-axis represents specific periods over stream. For example, P1 is the period from beginning to the first retraining, and P2 is the the period between first and second retraining etc.

Figure 6.2 on the next page shows the run time used for both stationary models and adaptive models for each data set. The updating takes more time for HTTP, SMTP+HTTP and ForestCover than PowerDemand and SMTP is due to that larger datasets take more time for prediction, and trigger more updating events. And for each retraining, the retrain buffers also contain larger retraining sets.

6.2 Synthetic example

In order to show the benefit of model updating over the stream, we demonstrate the on-line learning process of the small set PowerDemand in this section. The PowerDemand dataset does not contain clear incremental or sudden concept drift, but the normal pat-

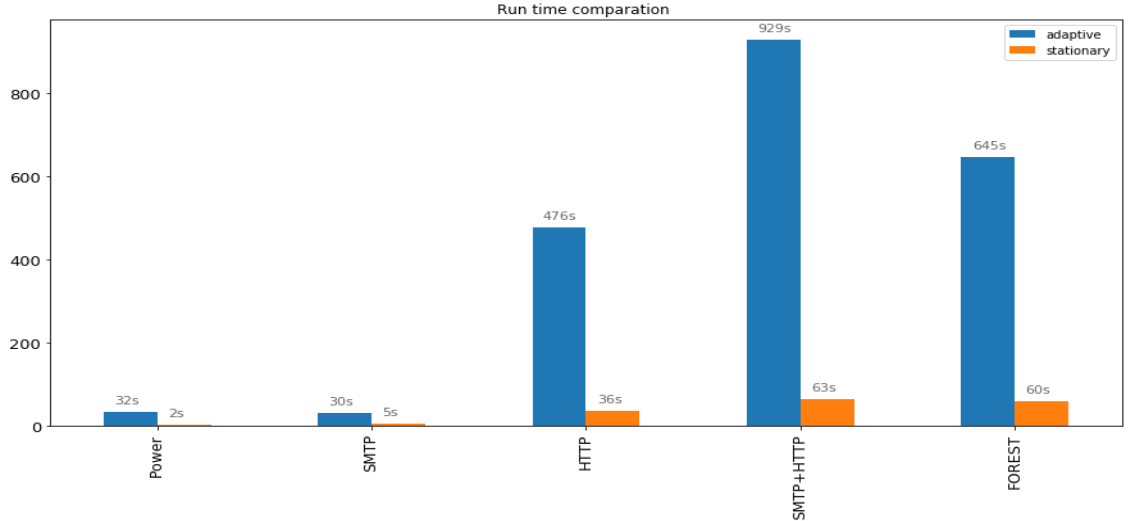


Figure 6.2: Run time statistics

terns still different slightly to each other. The lack of overall impression during the model initialization phase can lead to failures during the online phase.

6.2.1 Reaction of concept drift

Figure 6.3 shows 3 continual days power demand in normal state. Due to the lack of knowledge for current pattern, the autoencoder reconstructs the input time series higher than desired on day 1 (left diagram). This could be caused by seasonal changes on the power demand, which is slightly, gradually, and would potentially cause misclassification. The increase of normal data reconstruction error makes the margin between two classification classes smaller, and harder to make decision. As a consequence, the model retraining process is triggered after the second day with last seen data in the buffers, and the model performs well again on the third day.

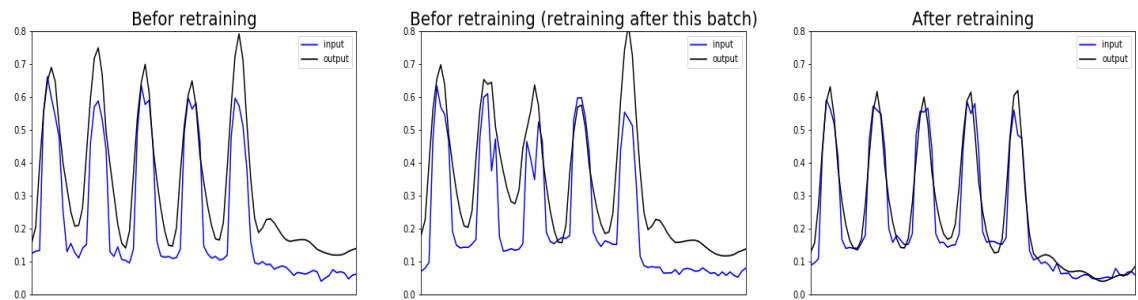


Figure 6.3: Retraining effect on PowerDemand dataset

6.2.2 Model updating

During the online phase, the model is retrained two times, before batch No.10 and No. 27. After retraining, the normal data reconstruction error becomes lower while for abnormal data becomes higher, so that the classification becomes easier.

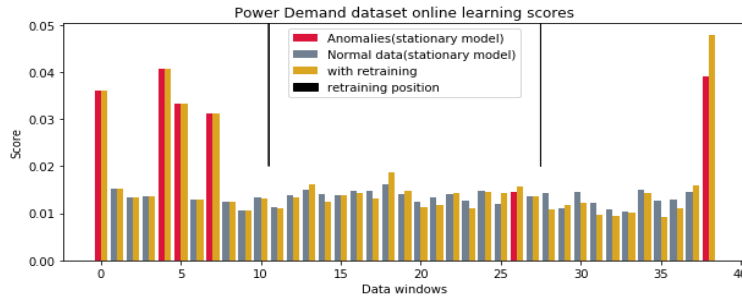


Figure 6.4: PowerDemand dataset online learning scores. For each window, the anomaly score is the highest pointwise scores within the window

After each updating process, the parameters μ , σ and threshold of anomaly scores are also updated. Figure 6.5 shows the parameter changes over the stream. As there is no clear concept drift during the power demand stream, the parameters change just slightly in order to learn latest knowledge from the retrain buffer.

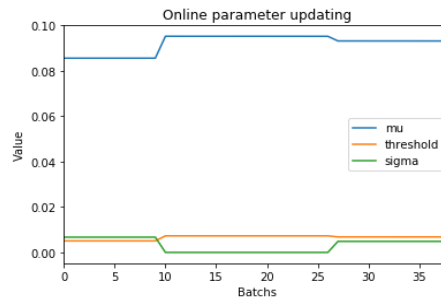


Figure 6.5: PowerDemand dataset online parameter updating

6.3 Model retraining

6.3.1 Reaction of sudden and drastic concept drift

The main advantage of online model is its ability to take reaction against sudden data distributional changes over time. The SMTP+HTTP data set is composed by directly connecting HTTP set after SMTP, so there is a sudden concept drift in between. The model is initialized with only SMTP data, so HTTP is completely unknown knowledge for the model. Figure 6.6 on the facing page shows the scores for both normal and abnormal data over the SMTP+HTTP stream. In the beginning, only SMTP data in the stream, and partially used for model initialization. In the online prediction phase, once the HTTP data arrives, the first peak of normal data scores' curve appears, and then the model updating is triggered, with buffer data being few hard SMTP data and most HTTP data. After the first model updating, the performance of model is still suboptimal due to the lack of enough HTTP data, therefore there are two further model updating process triggered during the following stream. As a result, the overall anomaly detection for SMTP+HTTP stream is good only except the short period after concept drift. The model updating are triggered in time after concept drift, and afterwards no redundant updating are triggered.

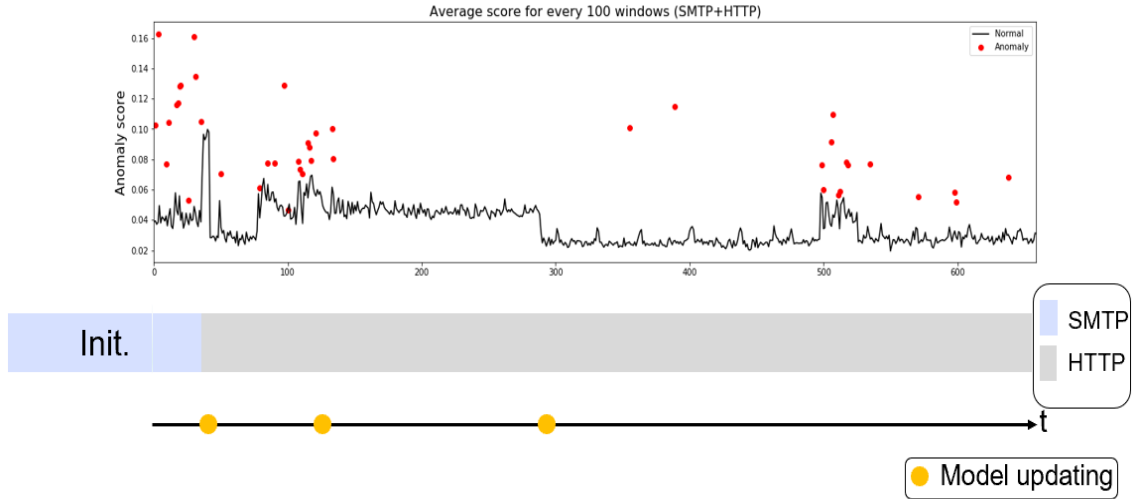


Figure 6.6: SMTP+HTTP stream concept drift.

6.3.2 Reaction of serial concept drift

Sometimes concept drift over the stream are slight, periodically, and potentially repeated. A single slight concept drift may not be able to trigger the retraining, but new knowledge should be saved into retraining buffer, so that once the model retrained with the fresh knowledge, the model should perform well when the same concept drift happens. We experiment with the ForestCover dataset. There are 7 kinds of forest cover types as labels. We take the least TYPE4 as anomaly while the rest 6 kinds as normal. The ForestCover stream is generated type by type, as shown in the bottom chart of Figure 6.7 on the next page. During the beginning phase, TYPE1 data appears in the stream, part of which is used for model initialization. Afterwards follows instances from TYPE2, TYPE3, TYPE5, TYPE6, TYPE7 and finally TYPE2 appears again during the ending phase. Anomaly data (TYPE4) is randomly distributed in the stream. Because the model is only initialized with TYPE1 data, every appearance of a new cover type will potentially cause a performance decrease. The concept drift under this setting is then the type changes over stream.

The points on the time axis in Figure 6.7 on the following page indicates the model updating. In the anomaly score chart, the scores for anomaly data are generally larger than normal data except when concept drifts take place. Once data stream from a new cover type appears in the stream, there are always peaks in the normal data score plot, and the difference to anomaly data scores decreases. After performance being impacted, the normal buffer is filled with hard windows shortly, that triggers the model updating quickly after the concept drift.

During the experiment, there are two cases delay the updating. Firstly, because of the fixed size of buffer, if a model updating is just triggered shortly before a concept drift, then the hard windows from new cover type need more time to fill the buffer and trigger updating. In Figure 6.7 on the next page, before TYPE3 arrive, there was a updating at the end

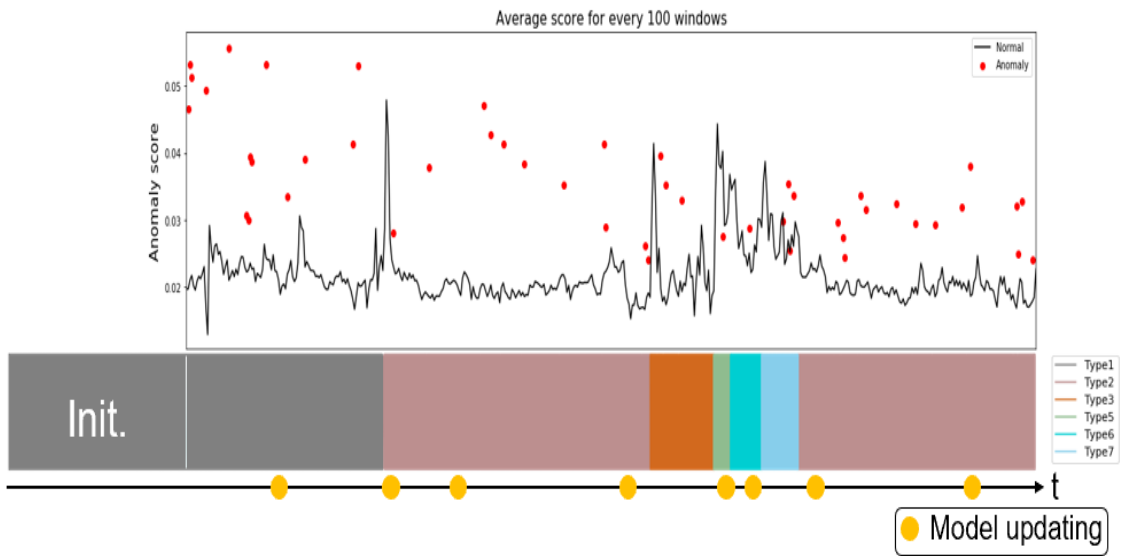


Figure 6.7: ForestCover stream concept drift. The chart on the bottom shows the cover type over the stream. The top chart shows average anomaly scores for every 100 windows

of TYPE2, and buffered was emptied, so that the model didn't take any action against the concept drift. Secondly, if a concept drift only appears in a short period, e.g. the TYPE7, which is also not enough to fill the buffer and trigger updating. However, under both aforementioned cases, the new information of concept drift, namely the new cover types, are stored in the buffer, and will be used for next model updating. If the concept drift missed model updating due to too short appearance period, we suppose that this would also not cause catastrophic effect over the stream prediction.

Chapter 7

Conclusion

Anomaly detection attracts more and more attention in the data mining field and have been applied to plenty of industrial use cases, which achieved perfect effectiveness and avoids large amount of financial spending. At the same time, the industrial applications need critically anomaly detection models under the big data background, specifically, ability to deal with high-volume, high-velocity data. In this paper, we proposed an adaptive LSTMs-autoencoder for streaming data anomaly detection. In the previous works, autoencoders are widely used in NLP tasks, e.g. language translation, sentence understanding. Vanilla autoencoders and deep autoencoders are also have been used to anomaly detection based on reconstruction error. [8] is the first work that use LSTMs-autoencoder for anomaly detection, with concentration to protection of temporal dependency between time series data. Our work uses similar LSTMs-autoencoder architecture, and enable the model to work with streaming data, and update model according to specific criterions. Our model shows good performance in detecting anomalies and outperforms the stationary models.

In terms of streaming data anomaly detection, we mainly focus on the concept drift over steam and model reinforcement by the last seen data. In the experiment with SMTP+HTTP dataset, our model shows robustness against sudden concept drift and adjusted the new data distribution very quickly. In the experiment with ForestCover dataset, the model masters serried and slight concept drifts also well. We also demonstrated an intuitive model online learning process with the small Power Demand dataset, which shows the impact of model updating between data windows in this small univariate dataset clearly.

Our model is designed under the assumption that there are expert labeling available during the online phase, which make the hard window collection become possible, and they are used for model updating. In the future work, a further research direction is to scale the model into fully automated without expert labeling online. Similar verification step as in [3] could be added after online prediction to make the model prediction more reliable, so that the data labeling can be directly according to the model prediction.

Bibliography

- [1] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder–decoder for statistical machine translation. 2014.
- [2] Hongyin Cui. Online outlier detection over data streams. 2002.
- [3] Yue Dong and Nathalie Japkowicz. Threaded ensembles of autoencoders for stream learning. 2017.
- [4] Adel Ghazikhani, Reza Monsefi, and Hadi Sadoghi Yadi. Online neural network model for non-stationary and imbalanced data stream classification. 2013.
- [5] Michiel Hermans and Benjamin Schrauwen. Training and analyzing deep recurrent neural networks. 2013.
- [6] Max Kochurov, Timur Garipov, Dmitry Podoprikin, Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. Bayesian incremental learning for deep neural networks.
- [7] Jiwei Li, Minh-Thang Luong, and Dan Jurafsky. A hierarchical neural autoencoder for paragraphs and documents. 2015.
- [8] Pankaj Malhotra, Anusha Ramakrishnan, Gaurangi Anand, Lovekesh Vig, Puneet Agarwal, and Gautam Shroff. Lstm-based encoder-decoder for multi-sensor anomaly detection. 2016.
- [9] Pankaj Malhotra, Vishnu TV, Lovekesh Vig, Puneet Agarwal, and Gautam Shroff. Timenet: Pre-trained deep recurrent neural network for time series classification. 2017.
- [10] Pankaj Malhotra, Lovekesh Vig, Gautam Shroff, and Puneet Agarwal. Long short term memory networks for anomaly detection in time series. 2015.
- [11] Marco Martinelli, Enrico Tronci, Giovanni Dipoppa, and Claudio Balducci². Electric power system anomaly detection using neural networks. 2004.
- [12] Mayu Sakurada and Takehisa Yairi. Anomaly detection using autoencoders with nonlinear dimensionality reduction. 2014.

- [13] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. 2014.
- [14] Swee Chuan Tan, Kai Ming Ting, and Tony Fei Liu. Fast anomaly detection for streaming data. 2011.
- [15] Xiaohong Tang and Chen Li. The stream detection based on local outlier factor. 2015.
- [16] Guanyu Zhou, Kihyuk Sohn, and Honglak Lee. Online incremental feature learning with denoising autoencoder. 2012.