

# CONTINUOUS MULTI-STEP PREDICTIONS OF HIGHLY IMBALANCED MULTIVARIATE TIME SERIES VIA DEEP LEARNING NETWORK

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1 Multi-step prediction of multivariate time series has always been a very popular  
2 research topic across industries. We focus on the scenario in which the data  
3 with severe imbalance problem caused by the 0 expansion in regression analysis,  
4 and meanwhile the data contains complex textual information. Such data is  
5 very common in customer’s life time value evaluation tasks in businesses. The  
6 commonly used two-stage modeling scheme effectively predicts whether or not a  
7 customer will pay for a product or service at the next moment. However, it is incapable  
8 of continuously forecasting potential payment values due to the strong imbalanced  
9 and randomness distribution of the data. In this paper, we propose a feature  
10 learning based deep learning method for imbalanced multivariate time series  
11 (FLIMTS). The innovative use of a weighted quantile loss in our proposed method  
12 handles the highly imbalance problem in regression. Furthermore, FLIMTS incorporates  
13 both the customer’s payment sequence and the behavioral characteristics of their  
14 interests which allows for more accurate predictions. Empirical analysis shows  
15 that FLIMTS has significant advantages and performs better than the existing  
16 two-stage approaches on common model evaluation criteria.

17 **Keywords:** highly imbalanced data, multivariate time series, LTV study, feature learning

## 18 1 INTRODUCTION

19 The Multi-step forecasting for multivariate time series has been widely used in daily life,  
20 such as finance, medicine, meteorology and, and other fields with very great commercial values.  
21 With the rapid development of computer science, data structure changes significantly. Multivariate  
22 time series data is usually accompanied with severe data imbalance problem of 0 expansions,  
23 which makes the multi-step predictions extremely difficult. We are interested in solving such  
24 multi-step ahead prediction problem of the informative multivariate time series data with  
25 severe imbalance problems in the data structure.

26  
27 The motivation of our research is from the statistical modeling and multi-step prediction of the  
28 customer’s life time value (LTV) sequence in commercial activities. We aim to study the impact  
29 of the user’s payment habits on the payment amount. The customer’s LTV refers to the total  
30 profit made from users during the time from product query to terminating the online services.  
31 As an important business indicator, the predicting the user’s payment values is a critical  
32 business requirement. Predicting the user’s payment directly determines the service providers’  
33 revenue capacity and their online service quality. The user’s payment value data is a  
34 typical extremely imbalanced multivariate time series. The highly imbalanced distribution  
35 of the response variable  $y_t$ , which refers to the user’s payment label value at the time  $t$ ,  
36 brings extreme difficulties to provide continuously multi-step ahead predictions of the  
37 multivariate time series data. This problem has been an unsolved challenge in industry  
for a long time.

38 In our scenario, the LTV prediction is a typical imbalanced regression analysis problem  
39 since the user’s payment value follows a highly imbalanced distribution with zero inflated.  
40 The proposed method should be conducted with the imbalanced learning technique. However,  
the available solu-

tions for the imbalanced problem are mostly designed for classification purposes. Even though Yang et al. (2021) has made some breakthroughs, the imbalance problem in regression is still a challenge question, and only a few achievements have been made in the related fields so far. Basically, there are two approaches to alleviate the imbalanced problem for classification, the data-based method and the model-based method. For the data-based method, undersampling in the majority class (Chawla et al., 2002; Han et al., 2005; He et al., 2008; Douzas and Bacao, 2019), and oversampling in the minority class (Chawla et al., 2002; Han et al., 2005; He et al., 2008; Douzas and Bacao, 2019) were used. However, these methods are not applicable to the regression problem directly, since the resampling method will bring strongly inductive bias to the distribution of the continuous label values. Compared to the data-based methods, the model-based methods might be applicable to the regression problems. For example, Lin et al. (2017); Li et al. (2019; 2020) added weights to samples or adjusted the objective (loss) function of the model. Yin et al. (2019); Huang et al. (2016); Yang and Xu (2020); Shu et al. (2019) used methods such as transfer learning, metric learning, and meta-learning techniques. Currently, the most competitive method for the imbalanced learning is by Kang et al. (2019). They decoupled the imbalanced learning into two stages of normal sampling in the feature learning stage, and balanced sampling in the Label Learning stage. The decoupled strategy achieves optimal modeling results so far.

Since the prediction of user’s payment values involves both regression and classification, the most commonly used solution in the industry is a two-stage approach, where the prediction process is disassembled into two sub-tasks: the classification task (stage-I: whether the user pays) and the regression task (stage-II: the payment amount of paying users) (Vanderveld et al., 2016; Chamberlain et al., 2017; Wang et al., 2019). Machine learning algorithms are adopted in these two tasks. Suppose the prediction result of the classification task is  $\hat{p}_i$ , and the given threshold is  $\omega$ , and the prediction result of the regression task is  $\hat{v}_i$ , then the model of the predicted paid value is  $\hat{y}_i = I(\hat{p}_i > \omega) \hat{v}_i$ . In the industry, the LightGBM algorithm (Ke et al., 2017) is commonly adopted for engineering implementation. We name this type of method as 2Stage-LGBM algorithm. The downside of these two-stage algorithms is that the model only can provide one time step ahead prediction due to the imbalance problem of the data. The current two-stage approaches achieve the multi-step ahead predictions in a clumsy way of that a series of independent predictive models for each target moment are established. This independent modeling scheme is not only difficult to conduct, but also a waste of time and computational resources. It increases the model maintenance cost with a unsatisfactory prediction accuracy. Therefore, we need a more sufficient predictive modeling strategy which allows for continuously multi-step ahead predictions in one model. Meanwhile, since we are analyzing the sequential data, Time should be included into the feature structure of the desired model.

In this article, we propose a deep learning algorithm based on the Feature Learning for imbalanced multivariate time series (FLIMTS). FLIMTS includes two parts: Representation Learning and Label Learning. Compared with the commonly used algorithms in industry, FLIMTS has two major advantages. First, by innovatively introducing a weighted quantile loss, it successfully eliminates the impact of the imbalanced data distributions. Second, FLIMTS fully utilizes the user’s portrait features, and deeply analyzes the static features and sequence features of the multivariate sequence data via feature learning processes to obtain more comprehensive sequence feature information. In the empirical analysis, we compare FLIMTS with the 2Stage-LGBM approach on a public data set. The results show that the proposed method performs better than the 2Stage-LGBM algorithm on common model evaluation criteria.

The rest of the article is organized as follows. Section 2 introduces the model and the details of the proposed algorithm. Section 3 is the empirical analysis. We compare our proposed algorithm with the 2Stage-LGBM algorithm on two data sets. Discussion in Section 4 concludes the article.

## 2 METHODOLOGY

### 2.1 MODEL AND NOTATION

The proposed method is a multi-step forecasting deep learning algorithm based on the feature learning for multivariate time series with heavy imbalance problem. Suppose that at time  $t$ ,  $\mathbf{y}_t^{(q)}$  are independent  $q$  step observations generated from the following imbalanced multivariate time series model:

$$\mathbf{y}_t^{(q)} = \mathcal{F}(\mathbf{x}_{t-p}, \dots, \mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{y}_t^{(p)} | \Theta) + \epsilon_t^{(q)}, \quad (1)$$

94 where  $\epsilon_t^{(q)}$  is the random error,  $\mathcal{F}$  is a unknown nonlinear mapping function, and  $\Theta = \{\Theta_{rep}, \Theta_{lab}\}$   
 95 is the parameter set of the Representation Learning module and the Label Learning module.  
 96  $p + 1$  is the size of the window of the previous data series used for later data series pre-  
 97 diction, and  $q$  is the length of the predicted series in the multi-step prediction.  $\mathbf{y}_t^{(p)} =$   
 98  $[y_{t-p} \ y_{t-p+1} \ \dots \ y_t]^T \in \mathbb{R}^{p+1}$  is the user paid value label sequence from time  $t - p$  to  
 99  $t$ .  $\mathbf{y}_t^{(q)} = [y_{t+1} \ y_{t+2} \ \dots \ y_{t+q}]^T \in \mathbb{R}^q$  is the user paid value label sequence from time  $t + 1$   
 100 to  $t + q$ .  $\mathbf{x}_t \in \mathbb{R}^m$  is the  $m$ -dimensional independent variable related to the feature variables at  
 101 time  $t$ . These features can be grouped into two feature vectors, which are the static feature vector  
 102  $\mathbf{x}_t^s \in \mathbb{R}^{m_1}$  and sequence feature vector  $\mathbf{x}_t^h \in \mathbb{R}^{m_2}$ , where  $m = m_1 + m_2$ . The static features do  
 103 not change over time. For example, some portrait features, such as gender and age in  $\mathbf{x}_t$ , can be  
 104 regarded as static features, where

$$\mathbf{x}_t^s = \mathbf{x}_{t'}^s = \mathbf{x}^s = [x_1^s \ x_2^s \ \dots \ x_{m_1}^s]^T, \quad \forall t \neq t'.$$

105 Observations of the sequence features change over time. For example, the number of user logins in  
 106 and the user's historical payment information  $\mathbf{y}_t^{(p)}$  are both sequence features, where

$$\mathbf{x}_t^h = [x_{t,1}^h \ x_{t,2}^h \ \dots \ x_{t,m'}^h \ y_{t-p} \ y_{t-p+1} \ \dots \ y_t]^T = [x_{t,1}^h \ x_{t,2}^h \ \dots \ x_{t,m_2}^h]^T,$$

107 where  $m_2 = m' + p + 1$ . Finally  $\mathbf{x}_t$  converts to

$$\mathbf{x}_t = \begin{bmatrix} \mathbf{x}_t^s \\ \mathbf{x}_t^h \end{bmatrix} = [x_1^s \ x_2^s \ \dots \ x_{m_1}^s \ x_{t,1}^h \ x_{t,2}^h \ \dots \ x_{t,m_2}^h]^T.$$

108 According to equation 1 we will use the sequence state information (including the static feature  
 109 information, the sequence feature information, and the label information) of the first  $p$  moments of  
 110 the data sequence to predict the label information of the following  $q$  moments of the data sequence.

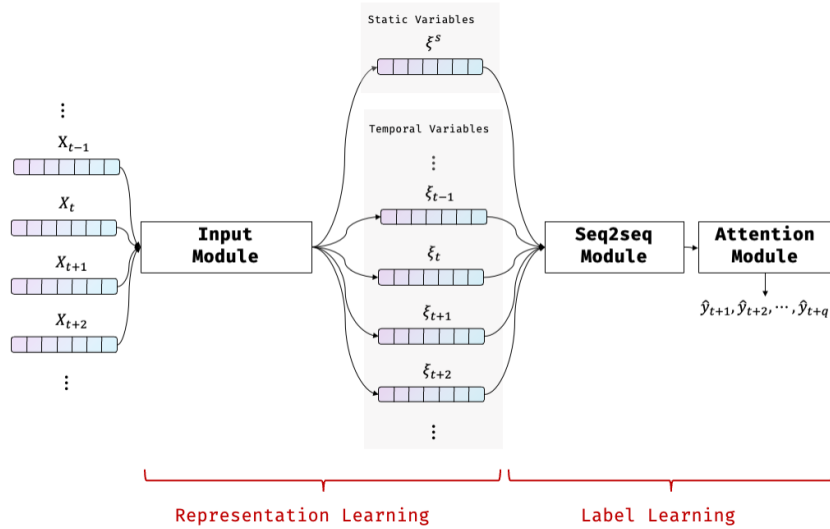


Figure 1: Model Architecture Diagram of FLIMTS

111 Figure 1 is the computational framework of the FLIMTS algorithm. The key idea of the FLIMTS  
 112 algorithm is to fit a nonlinear mapping  $\mathcal{F}$  using a deep neural network. The algorithm structure  
 113 includes two parts: the Representation Learning module and the Label Learning Module. In the

114 Representation Learning part, the deep neural network combines continuous variables and discrete  
 115 variables to achieve a desired analysis effect. The Representation Learning part is to decompose  
 116 the independent variable sequence into static features and sequence features, and then convert them  
 117 into the continuous representation vectors. The Label Learning part includes the Seq2seq module,  
 118 the Attention module and the the Output module. It maps the continuous representation vectors  
 119 generated from the Representation Learning to the target value of the response variables. For each  
 120 sequence data sample, it will be converted into a vector through the Representation Learning module  
 121 firstly, then will be processed by the Seq2seq module, the Attention module and the output module  
 122 to obtain the Multi-step prediction values. The computational workflow of the proposed algorithm  
 123 FLIMTS is summarized in Algorithm 1. The details of our designed deep learning networks are  
 124 introduced in Appendix A.

---

**Algorithm 1: FLIMTS**


---

**Input:**  $\{\mathbf{x}_\tau\}$ , for  $\tau \in [t-p, t]$

**Output:** the  $q$ -step predictions of the sequence labels  $\hat{\mathbf{y}}_t^{(q)} = [\hat{y}_{t+1}, \hat{y}_{t+2}, \dots, \hat{y}_{t+q}]^T$

**Step 1 Representation Learning:**

(a) Obtain the embedded static feature vector  $\xi^s$  of the input variable through the static feature processing module (Algorithm 2 in Appendix A.1):

$$\xi^s \leftarrow StaInput(\mathbf{x}^s);$$

(b) Obtain the embedded sequence feature vector  $\xi_\tau^h$  of the input variable through the sequence feature processing module (Algorithm 3 in Appendix A.1):

$$\xi_\tau^h \leftarrow TempInput(\mathbf{x}_\tau^h), \quad t-p \leq \tau \leq t.$$

**Step 2 Label Learning:**

**Step 2.1** process the feature vector through the Seq2seq module (Algorithm 4 in Appendix A.2.1):

$$[\mathbf{h}_{t-p} \dots \mathbf{h}_t \dots \mathbf{h}_{t+q}]^T \leftarrow Seq2seqModule\left([\xi^s, \xi_{t-p}^h, \dots, \xi_{t-1}^h, \xi_t^h]^T\right);$$

**Step 2.2** predict the  $q$ -step sequence labels vector  $\hat{\mathbf{y}}_t^{(q)}$  through the Attention module (Algorithm 5 in Appendix A.2.2):

$$\hat{\mathbf{y}}_t^{(q)} \leftarrow AttentionModule\left([\mathbf{h}_{t-p} \dots \mathbf{h}_t \dots \mathbf{h}_{t+q}]^T\right).$$


---

125 2.2 WEIGHTED QUANTILE LOSS AND MODEL OPTIMIZATION

126 The proposed deep learning network contains a large number of parameters. We adopt a two-stage  
 127 parameter optimization strategy for the proposed method, where the parameters of the Representa-  
 128 tion Learning model and the Label Learning model are optimized separately and sequentially. In the  
 129 Representation Learning stage, we use the common Quantile Loss function. In the Label Learning  
 130 stage, to deal with the highly imbalanced problems of the multivariate time series, we innovatively  
 131 introduce a Weighted Quantile Loss function, which greatly improves the effect of multi-step fore-  
 132 casting. Then we iteratively update all parameters based on the gradient back propagation technique.

133 The Representation Learning is the first stage of the proposed FLIMTS algorithm, since we do not  
 134 have to consider the prior information of the label distribution, we only need to learn the pattern of  
 135 the data distribution. Therefore, we do not have to do subsampling in this modeling stage. Balanced  
 136 subsampling is enough if necessary. In the learning process, the optimal estimates of all parameters  
 137 are obtained by minimizing the following Quantile Loss

$$QLoss(\eta, y_i, \hat{y}_i) = \max\{\eta \cdot (\hat{y}_i - y_i), (1 - \eta) \cdot (y_i - \hat{y}_i)\},$$

138 where  $\eta \in (0, 1)$  adjusts the prediction tendency of the algorithm. If  $\eta$  is large, the algorithm  
 139 parameters will be trained in the direction of underestimating the label value. If  $\eta$  is small, it will  
 140 be trained in the direction of overestimating the label value. When  $\eta$  is set to 0.5, the effect of the  
 141 Quantile Loss will be equivalent to the Absolute Value Loss.

142 The Label Learning phase is the second stage of the proposed FLIMTS algorithm. The purpose  
 143 of this stage is to fit the mapping relationship between the feature vectors and the label val-  
 144 ues. Therefore, the prior knowledge of the label distribution will have a significant impact on  
 145 the modeling process. In the Label Learning phase, parameters are optimized by minimizing the  
 146 Weighted Quantile Loss. To construct the Weighted Quantile Loss, we need to obtain its asymp-  
 147 totic distribution of the label value. Suppose that the partition of the range of the label value is  
 148  $x_{(0)} < x_{(2)} < \dots < x_{(j)} < \dots < x_{(N)} = \infty$ . We estimate the distribution of the label value  $F(y)$   
 149 by the empirical distribution function  $G_N(y) = \frac{1}{N} \sum_{j=1}^N I(Y_j \leq y)$ . The weight of the  $i^{th}$  label  
 150 value is

$$w_i = \sum_{j=0}^N \frac{I(x_j \leq y_i < x_{j+1})}{\int_{x_j}^{x_{j+1}} f(x) dx}, \quad i = 1, \dots, N.$$

151 Then the Weighted Quantile Loss is

$$WQLoss(\eta, Y, \hat{Y}) = \sum_{i=1}^q w_i \cdot QLoss(\eta, y_i, \hat{y}_i).$$

152 The Weighted Quantile Loss function is designed to be sensitive to the skewness of the target label  
 153 distribution and the sensitivity is reinforced for intervals with few samples, since the total loss com-  
 154 ing from these parts are highly likely to be underrated due to its small quantity of samples in the  
 155 optimization process. In reality, in order to improve programming efficiency and reduce the time  
 156 expense caused by the memory exchange in the communication between different devices (CPU  
 157 and GPU), our empirical distribution function  $G_N(y)$  is usually obtained by integrating the results  
 158 of each batch in the parallel computing process. This method may cause small amount of deviation  
 159 when estimating the distribution of the data with imbalance problem. However, this small deviation  
 160 will gradually converge as the number of training rounds and the sample size increases. Therefore  
 161 it will have little impact on the overall fittings of the model.

162 Let  $\Theta = \{\Theta_{rep}, \Theta_{lab}\}$  is the set of parameters that need to be optimized in the proposed model,  
 163 where  $\Theta_{rep}$  is the parameter set of the Representation Learning part,  $\Theta_{lab}$  is the parameter set of the  
 164 Label Learning part. These two sets of the parameters cannot be completely separated in the training  
 165 process. In the Representation Learning part,  $\Theta_{lab}$  is also updated while optimizing  $\Theta_{rep}$ . Mean  
 166 while, in order to improve the stability of parameter estimation, we need to optimize the collective  
 167 loss under different values of  $\eta$ . In our empirical analysis we use  $\eta \in \{0.3, 0.5, 0.7\}$ . Let  $L_{rep}$  be  
 168 the Quantile Loss of the Representation Learning.

$$L_{rep} = \sum_{\eta} \sum_{i=1}^q QLoss(\eta, \hat{y}_{t+i}, y_{t+i}),$$

169  $\Theta_{rep}$  is optimized and updated by the minimizer of  $L_{rep}$ .

$$\Theta_{rep} = \Theta_{rep}^* - \alpha \frac{\partial L_{rep}}{\partial \Theta_{rep}^*}, \quad \Theta_{lab(rep)} = \Theta_{lab(rep)}^* - \alpha \frac{\partial L_{rep}}{\partial \Theta_{lab(rep)}^*},$$

170 where  $\Theta^*$  is the previous state of the parameters, and  $\Theta$  is the updated state of the parameters.  
 171  $\Theta_{lab(rep)}$  is the parameter set of the Label Learning updated in the Representation Learning stage.  
 172 Let  $L_{lab}$  be the Weighted Quantile Loss of the Label Learning.

$$L_{lab} = \sum_{\eta} \sum_{i=1}^q WQLoss(\eta, \hat{y}_{t+i}, y_{t+i}).$$

173  $\Theta_{lab}$  is firstly initialized by  $\Theta_{lab(rep)}$ , that is  $\Theta_{lab}^{(0)} = \Theta_{lab(rep)}$ . Then  $\Theta_{lab}$  is optimized and updated  
174 by

$$\Theta_{lab} = \Theta_{lab}^* - \alpha \frac{\partial L_{lab}}{\partial \Theta_{lab}^*}.$$

175 In addition, the generalization performance of the model is improved by controlling the Dropout  
176 Rate (*DPR*).

### 177 3 EMPIRICAL ANALYSIS

178 We compare the proposed method with the commonly used 2Stage-LGBM algorithm on the public  
179 dataset AVSC. The 2Stage-LGBM algorithm is implemented based on the LightGBM algorithm  
180 and is one of the best prediction scheme in the industry for multivariate time series with highly  
181 imbalanced problems. Notice that, in our experiment we make four step forward predictions from  
182 time  $t + 1$  to  $t + 4$ . FLIMITS only needs to build a single model to generate continuous multi-step  
183 predictions, while the 2Stage-LGBM has to build four separate prediction models for each target  
184 moment.

185 The AVSC dataset is a desensitized transaction data (Dua and Graff, 2017) from the Acquire Valued  
186 Shoppers Challenge. This dataset includes transactions of some brick-and-mortar stores over a  
187 period of time with 11 features for each data record, such as User ID, Store ID, Item Category, Item  
188 Subcategory, Company ID, Brand ID, Purchase Date, Number of Purchased Item, Measurement  
189 Unit of the Purchased Item, and Purchase Amount. The original data structure of the AVSC dataset  
190 is not very suitable for the purpose of our analysis since it does not show out the user’s consumption  
191 behaviors very well. The data must be preprocessed before performing any further analysis. Let’s  
192 define the user’s payment value is the amount of repurchase made by a user who has a purchase  
193 record before. We first clean and aggregate the original data based on the customer ID, and obtain  
194 the monthly customer’s consumption data, which is the overall paid value of each customer in 8  
195 months. The reprocessed data has two groups of features, the Payment features and the Context  
196 features. The Payment features are used to describe the customer’s payment behavior, including  
197 the number of payments, the frequency of payments, and the average amount of each payment, etc.  
198 The Context features are the statistical characteristics of existing customers in a store, such as the  
199 payment amount per capita and the payment frequency per capita, etc. Since this is a monthly data,  
200 the data at each time point is the cumulative purchase value within each month. Not surprisingly, the  
201 paid value of is extremely imbalanced, since the ratio of non-paying users to paying users exceeds  
202 3 : 1. Among these paying users, the number of customers with higher paid value decreases sharply  
203 as the paid value increases. To conform the data to the real business scenario, we split the dataset  
204 into the training set and the test set based on the dates. We use the data before Dec, 1<sup>st</sup>2012 as the  
205 training set, which has 1,355,880 data records. Then the test set size is 356,980. Their ratio is about  
206 8 : 2.

207 Since the user payment value is a continuous variable, we choose the mean square error (*MSE*)  
208 and the mean absolute error (*MAE*), which are two commonly used model evaluation criteria in  
209 regression. In reality, the downstream tasks of paid value prediction are usually related to the ranking  
210 of users (such as selecting TOP-N from users for advertising, etc.). Therefore we adopt the  $rAUC =$   
211  $P(\hat{y}_1 > \hat{y}_2 | y_1 > y_2)$  (*regression - AUC*) criterion, which measures the ranking quality of the  
212 regression model. A better model can rank samples with larger true label values ahead of samples  
213 with smaller true label values when sorting the samples according to their predicted values from  
214 large to small. The larger the *rAUC* is, the higher the accuracy of the algorithm in sorting users  
215 according to the predicted value, and vice versa.

216 The computational environment of experiments in this article are: *CPU*Interi7 –  
217 10700; *GPU*NvidiaRTX3060; *RAM*32Gb; *Software*Python3.9 + *CUDA*11.1. We use the

218 open source *SQL* query engines (Impala and Trino) for data cleaning and feature engineering. The  
 219 implementation of algorithm engineering relies on the Pytorch (Paszke et al., 2019) environment,  
 220 and we also uses the Pytorch-lightning module for parameter optimization.

221 3.1 EXPERIMENT RESULTS

222 In the experiment we set the parameters  $p = 6, q = 4$  in the model equation 1. The initialization  
 223 settings of hyperparameters are: the learning rate is 0.0001; the number of LSTM layer is 1; the  
 224 dimension of the intermediate and hidden variables in the Seq2seq module is  $d_k = 128$ ; the number  
 225 of heads in the Multi-head Self Attention module is  $H = 8$ ; the dropout rate is  $DPR = 0.2$ .

226 The results of the model evaluation criteria on the training set are shown in Figure 2. The trend of  
 227 all criteria drops rapidly in the early stage of training, and tends to be a stable fluctuation later. This  
 228 indicates that the proposed model converges fast in the training stage.

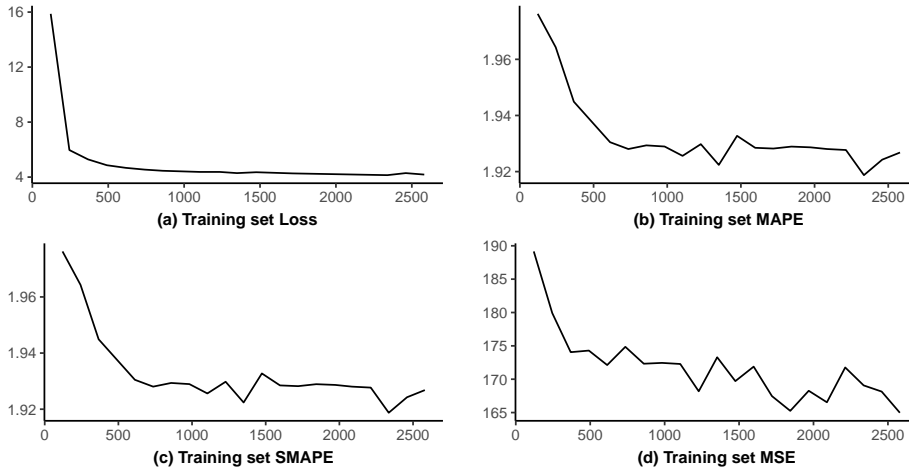


Figure 2: The convergence rate of the proposed model on the training stage.

229 In practice, the predictive quality of those paying users can better reflect the advantages of the  
 230 models, thus attract more attention in business. We compare the model performance of these two  
 231 algorithms in predicting the paid value at the  $(t + 1)th$  moment on two customer groups, which are  
 232 the All User group and the Paying User group. The result is summarized in the Table 1. We see that  
 233 the MSE and MAE of FLIMITS are much smaller than the 2Stage-LGBM approach, which means  
 234 that the mean and median of the predicted paid value by FLIMITS are closer to the true value than  
 235 the 2Stage-LGBM algorithm for both All User group and Paying User group. In terms of the rAUC,  
 236 since both model schemes are based on regression analysis, therefore the rAUC of FLIMITS is very  
 237 close to the 2Stage-LGBM algorithm, which is a reasonable result.

Table 1: Model Performance Comparison

	Algorithm	MSE	MAE	rAUC
All User	FLIMITS	<b>51.54</b>	<b>0.85</b>	<b>0.98</b>
	2Stage-LGBM	151.05	0.91	0.98
Paying User	FLIMITS	<b>208.82</b>	<b>1.55</b>	0.90
	2Stage-LGBM	615.32	2.18	<b>0.90</b>

238 Figure 3 shows the MSE, MAE and rAUC of the FLIMITS and the 2Stage-LGBM algorithms in  
 239 multi-step ahead predictions of the paid values at the next four moments of time  $t+1$  to  $t+4$ . In terms  
 240 of the rAUC, the proposed algorithm shows a slight advantage over the 2Stage-LGBM algorithm as  
 241 the size of the prediction time step increases. For the MSE and MAE, the proposed algorithm  
 242 provides much smaller results than that of the 2Stage-LGBM algorithm, and this advantage grows  
 243 as the size of the prediction time step increases.

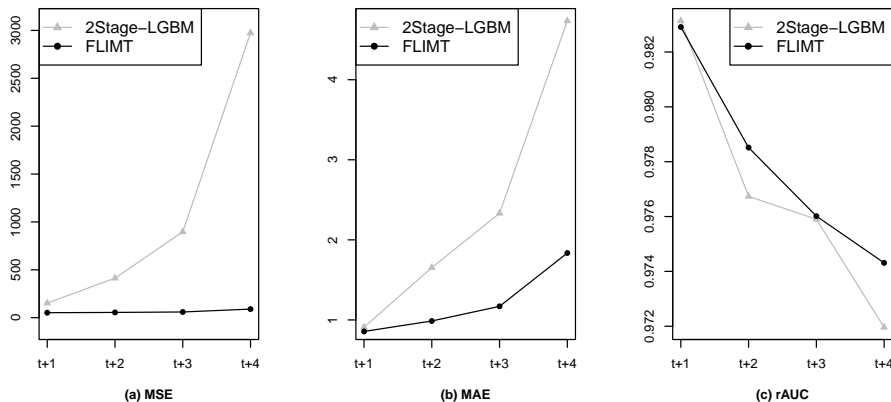


Figure 3: The MSE, MAE and rAUC of different algorithms in multi-step predictions of the paid value at the next four moments of time  $t + 1$  to  $t + 4$ .

244 Figure 4 shows the multi-step prediction results of the paid value from time  $t + 1$  to  $t + 4$  based on  
 245 the FLIMITS algorithm and the 2Stage-LGBM algorithm for four types of paying value users, which  
 246 are the High paying value users, Mid paying value users, Low paying value users, and Null paying  
 247 value users. The prediction results at all four moments show that FLIMITS has better predicting  
 248 performance than the 2Stage-LGBM method for all four types of paying value users with much less  
 249 cost of the computational resources, since FLIMITS only needs to build a single model to generate  
 250 multi-step predictions continuously, while the 2Stage-LGBM has to build four separate prediction  
 251 models for each target moment from  $t + 1$  to  $t + 4$ .

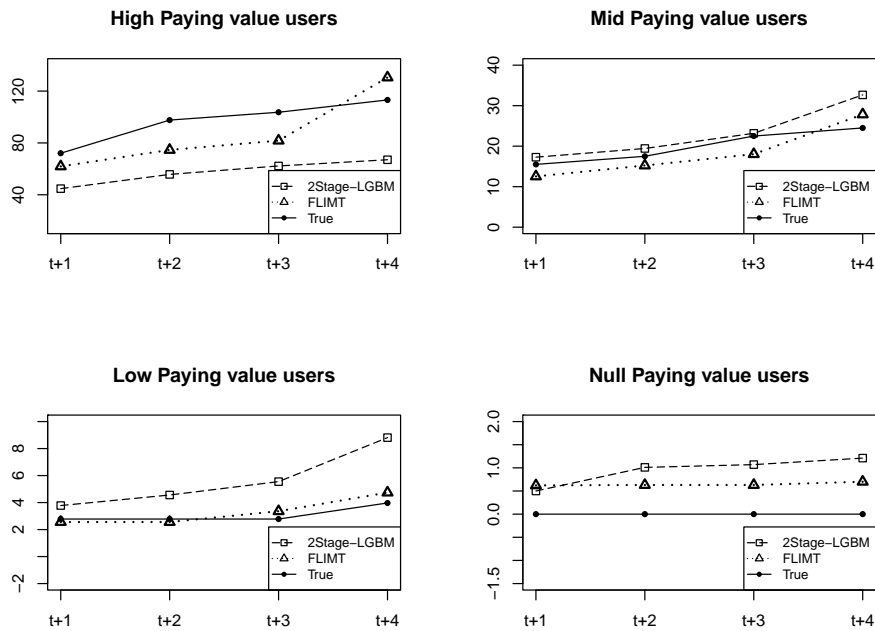


Figure 4: Multi-step predictions of the paid value at the next four moments from  $t + 1$  to  $t + 4$  for four types of paying value users under two methods.



252 4 CONCLUSION

253 In this paper we propose a new deep learning algorithm, FLIMTS, for the multi-step forecasting  
254 of the multivariate time series with severe data imbalance problem. The great advantage of the  
255 FLIMTS algorithm is that by introducing a weighted quantile loss, we greatly reduce the influence  
256 of the imbalanced data distribution problem. Therefore, the proposed method only needs to build  
257 one predictive model to generate multi-step ahead predictions for a sequence of target moments for  
258 imbalanced multivariate time series. In contrast, the traditional 2Stage-LGBM algorithm must build  
259 separate predictive models at each target moment to achieve satisfactory multi-step predictions. Ad-  
260 ditionally, the prediction accuracy is improved by the rigorously designed deep learning networks,  
261 which combine the Representation Learning and Label Learning based on the Feature Learning  
262 techniques.

263 AUTHOR CONTRIBUTIONS

264 ACKNOWLEDGMENTS

## 265 REFERENCES

- 266 Chamberlain, B. P., A. Cardoso, C. B. Liu, R. Pagliari, and M. P. Deisenroth (2017). Customer life-  
267 time value prediction using embeddings. In *Proceedings of the 23rd ACM SIGKDD international*  
268 *conference on knowledge discovery and data mining*, pp. 1753–1762.
- 269 Chawla, N. V., K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer (2002). Smote: synthetic minority  
270 over-sampling technique. *Journal of artificial intelligence research* 16, 321–357.
- 271 Douzas, G. and F. Bacao (2019). Geometric smote a geometrically enhanced drop-in replacement  
272 for smote. *Information sciences* 501, 118–135.
- 273 Dua, D. and C. Graff (2017). Acquire valued shoppers challenge.
- 274 Han, H., W.-Y. Wang, and B.-H. Mao (2005). Borderline-smote: a new over-sampling method in  
275 imbalanced data sets learning. In *International conference on intelligent computing*, pp. 878–887.  
276 Springer.
- 277 He, H., Y. Bai, E. A. Garcia, and S. Li (2008). Adasyn: Adaptive synthetic sampling approach  
278 for imbalanced learning. In *2008 IEEE international joint conference on neural networks (IEEE*  
279 *world congress on computational intelligence)*, pp. 1322–1328. IEEE.
- 280 He, K., X. Zhang, S. Ren, and J. Sun (2016). Deep residual learning for image recognition. In *2016*  
281 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778.
- 282 Huang, C., Y. Li, C. C. Loy, and X. Tang (2016). Learning deep representation for imbalanced  
283 classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
284 pp. 5375–5384.
- 285 Ioffe, S. and C. Szegedy (2015). Batch normalization: Accelerating deep network training by re-  
286 ducing internal covariate shift.
- 287 Kang, B., S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis (2019). Decoupling  
288 representation and classifier for long-tailed recognition. *arXiv preprint arXiv:1910.09217*.
- 289 Ke, G., Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu (2017). Lightgbm: A  
290 highly efficient gradient boosting decision tree. In *Proceedings of the 31st International Confer-*  
291 *ence on Neural Information Processing Systems, NIPS’17*, Red Hook, NY, USA, pp. 31493157.  
292 Curran Associates Inc.
- 293 Li, B., Y. Liu, and X. Wang (2019). Gradient harmonized single-stage detector. In *Proceedings of*  
294 *the AAAI Conference on Artificial Intelligence*, Volume 33, pp. 8577–8584.
- 295 Li, X., W. Wang, L. Wu, S. Chen, X. Hu, J. Li, J. Tang, and J. Yang (2020). Generalized focal loss:  
296 Learning qualified and distributed bounding boxes for dense object detection. *arXiv preprint*  
297 *arXiv:2006.04388*.
- 298 Lin, T., P. Goyal, R. B. Girshick, K. He, and P. Dollár (2017). Focal loss for dense object detection.  
299 *CoRR abs/1708.02002*.
- 300 Paszke, A., S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimeshein,  
301 L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy,  
302 B. Steiner, L. Fang, J. Bai, and S. Chintala (2019). Pytorch: An imperative style, high-  
303 performance deep learning library. In *Advances in Neural Information Processing Systems 32*,  
304 pp. 8024–8035. Curran Associates, Inc.
- 305 Shu, J., Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng (2019). Meta-weight-net: Learning an  
306 explicit mapping for sample weighting. *arXiv preprint arXiv:1902.07379*.
- 307 Vanderveld, A., A. Pandey, A. Han, and R. Parekh (2016). An engagement-based customer lifetime  
308 value system for e-commerce. In *Proceedings of the 22nd ACM SIGKDD international conference*  
309 *on knowledge discovery and data mining*, pp. 293–302.
- 310 Wang, X., T. Liu, and J. Miao (2019). A deep probabilistic model for customer lifetime value  
311 prediction. *arXiv preprint arXiv:1912.07753*.

- 312 Yang, Y. and Z. Xu (2020). Rethinking the value of labels for improving class-imbalanced learning.  
 313 *arXiv preprint arXiv:2006.07529*.
- 314 Yang, Y., K. Zha, Y.-C. Chen, H. Wang, and D. Katabi (2021). Delving into deep imbalanced  
 315 regression. *arXiv preprint arXiv:2102.09554*.
- 316 Yin, X., X. Yu, K. Sohn, X. Liu, and M. Chandraker (2019). Feature transfer learning for face recog-  
 317 nition with under-represented data. In *Proceedings of the IEEE/CVF Conference on Computer  
 318 Vision and Pattern Recognition*, pp. 5704–5713.

## 319 A APPENDIX

### 320 A.1 REPRESENTATION LEARNING NETWORK

321 The Representation Learning is designed to learn the information of all features and convert them  
 322 into the representation vectors for the further quantitative analysis. At each time  $t$ , the feature vector  
 323 of each sample is  $\mathbf{x}_t = (x_1^s, x_2^s, \dots, x_{m_1}^s, x_{t,1}^h, x_{t,2}^h, \dots, x_{t,m_2}^h)$ , where each  $x$  represents the obser-  
 324 vation of the corresponded feature at time  $t$ . The algorithm takes the sequence features as the input  
 325 to the Encoder. In the modeling process, since these features are not necessarily numerical continu-  
 326 ous, and discrete features cannot be directly used in the deep neural network model, it is necessary  
 327 to use the embedding method to convert these features into vectors with specified dimension. And  
 328 this is also the main purpose of using the Representation Learning. In the following content, we use  
 329  $\xi^s$  and  $\xi_t^h$  to represent the static feature vectors and the sequence feature vectors that are learned  
 330 from the Representation Learning respectively.

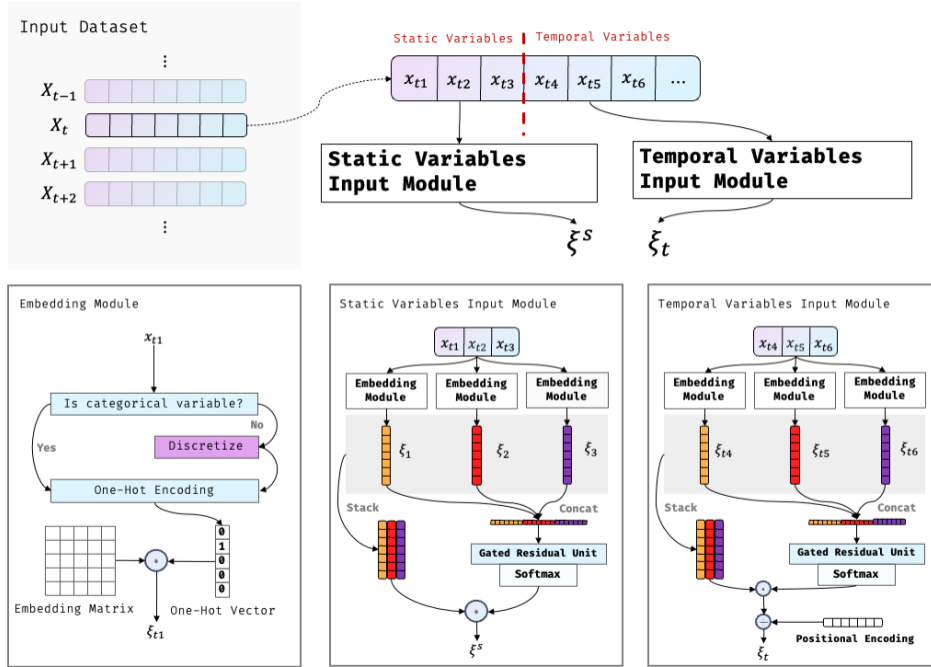


Figure 5: The operational flow chart of the Representation Learning

331 Suppose we set the embedding vector  $\xi_i$  at time  $t$  to be a  $d_k$  dimensional vector. For the  $i^{th}$  discrete  
 332 feature  $X_{t,i}$ , we encode the feature firstly, and then find the corresponded  $d_k$ -dimensional embedding  
 333 vector  $\xi_i$  from the Embedding Matrix of  $x_{t,i}$ , where  $d_k$  is a hyperparameter that has to be decided  
 334 in advance. For example, in the following empirical analysis section, we set  $d_k = 128$ . For the  
 335  $j^{th}$  continuous feature  $X_{t,j}$ , according to the specific situation, we could transform it into a discrete  
 336 variable using the binning technique or convert it into a  $d_k$  dimensional vector  $\xi_j$  through a  $d_k$

337 dimensional fully connected network layer. At this stage, both the parameters of the Embedding  
338 Matrix and the fully connected network can be optimized in the training process.

339 Figure 5 is the operational flow chart of the Representation Learning. At time  $t$ , the Representation  
340 Learning divides the input vector  $\mathbf{x}_t$  into the static features and sequence features, and processes  
341 them separately. For the static features, we designed a static feature processing algorithm, StaInput,  
342 and its computational workflow is described in the Algorithm 2. For the dynamic sequence features,  
343 we designed a sequence feature processing algorithm, TempInput, and its computational workflow  
344 is described in the Algorithm 3. Both algorithms map each feature into a dense embedding vector  
345 via the embedding process,

$$\begin{aligned}\xi_i^s &\leftarrow Emb_i(x_i^s), \\ \xi_{t,j}^h &\leftarrow Emb_i(x_{t,j}^h),\end{aligned}$$

346 where  $i \in [1, m_1]$ ,  $\xi_i^s \in \mathbb{R}^{d_k \times 1}$ , and  $j \in [1, m_2]$ ,  $\xi_{t,j}^h \in \mathbb{R}^{d_k \times 1}$ . Then obtain the representation  
347 vectors by the weighted average of these embedding vectors. The difference is the representation  
348 learning for the sequence features have an additional temporal information encoding step.

349 To more accurately extract the representation information of the sequence, the Representation Learn-  
350 ing estimates the weight of the embedding vector, and uses the weighted average method to calculate  
351 the representation vectors  $\xi^s$  and  $\xi_t^h$ . In this specific process, the linear transformation is firstly per-  
352 formed after splicing each vector vertically to convert it into a  $d_k$  dimensional vector. Then calculate  
353 the weights of the features based on the normalized vector using the Softmax module. If the weight  
354 of a feature is close to 0, it means the importance of this feature is very low. Otherwise its importance  
355 is very high. The representation vectors are then updated in the following way,

$$\begin{aligned}\xi^s &= [\xi_1^s \xi_2^s \dots \xi_{m_1}^s] Softmax\left(W_1 [\xi_1^{sT} \xi_2^{sT} \dots \xi_{m_1}^{sT}]^T + \mathbf{b}_1\right), \\ \xi_t^h &= [\xi_{t,1}^h \xi_{t,2}^h \dots \xi_{t,m_2}^h] Softmax\left(W_2 [\xi_{t,1}^{hT} \xi_{t,2}^{hT} \dots \xi_{t,m_2}^{hT}]^T + \mathbf{b}_2\right),\end{aligned}\tag{2}$$

356 where  $\xi^s \in \mathbb{R}^{d_k \times 1}$ ,  $W_1 \in \mathbb{R}^{m_1 \times (d_k \cdot m_1)}$ ,  $\mathbf{b}_1 \in \mathbb{R}^{m_1 \times 1}$ , and  $\xi_{t,m}^h \in \mathbb{R}^{d_k \times 1}$ ,  $W_2 \in \mathbb{R}^{m_2 \times (d_k \cdot m_2)}$ ,  
357  $\mathbf{b}_2 \in \mathbb{R}^{m_2 \times 1}$ .

358 When processing the sequence features, since the subsequent structure of the algorithm includes  
359 the Attention module, it is necessary to assign positional encoding to the input vector, that is, add  
360 the temporal information of the  $t^{th}$  moment to the sequence feature vector  $\xi_t^h$ . Suppose that the  
361 temporal information at the  $t^{th}$  moment is  $TimeSeq_t$ , where

$$TimeSeq_t = (t - p, t - p + 1, \dots, t, \dots, t + q).$$

362 In our application, the user life cycle sequence is a non-negative monotonic increasing sequence.  
363 For different products, users' payment habits show obvious personalized patterns. For example, the  
364 payment behavior is concentrated in the early or late stage of the customer's life cycle. Therefore,  
365 we hope to convert the  $t^{th}$  moment's time tag information  $TimeSeq_t$  into a representation vector  
366 of the feature sequence by the Embedding method, and then add it as a position code into  $\xi_t^h$ . This  
367 is the major difference of processing static features and sequential features. The popular positional  
368 encoding composes of sine and cosine functions, which is not suitable here. In our TempInput  
369 algorithm the positional encoding is obtained via model training. The specific method uses the  
370 Embedding technique to convert the temporal label sequence  $TimeSeq_t$  at the  $t^{th}$  moment into a  
371 discrete feature vector, and then superimposes it into the output vector  $\xi_t^h$ , that is

$$\xi_t^h = \tilde{\xi}_t^h + Emb(TimeSeq_t),$$

372 where  $\tilde{\xi}_t^h$  is the representation vector obtained from the previous step by equation 2. In addition, the  
373 position encoding also needs to be added on the Decoder side. The input of the Decoder module is  
374 denoted by  $\xi_t^f$  in the later sections.

**Algorithm 2: StaInput****Input:** The static feature data of user  $\mathbf{x}^s = (x_1^s, x_2^s, \dots, x_{m_1}^s)$ **Output:**  $\xi^s$  the embedding vector of static features**Step 1:** For  $i = 1 : m_1$ , compute embedding vectors:

$$\xi_i^s \leftarrow Emb_i(x_i^s)$$

**Step 2:** Take a weighted average of the embedding vectors:

$$\xi^s = [\xi_1^s \quad \xi_2^s \quad \dots \quad \xi_{m_1}^s] Softmax\left(W_1 [\xi_1^{sT} \quad \xi_2^{sT} \quad \dots \quad \xi_{m_1}^{sT}]^T + \mathbf{b}_1\right)$$

**Algorithm 3: TempInput****Input:** The sequence feature data of the user at time  $t$   $\mathbf{x}_t^h = (x_{t,1}^h, x_{t,2}^h, \dots, x_{t,m_2}^h)$ ,  
 $TimeSeq_t = (t - p, t - p + 1, \dots, t, \dots, t + q)$ **Output:**  $\xi_t^h$  the embedding vector of sequence features at time  $t$ **Step 1:** For  $j = 1 : m_2$ , computing embedding vectors:

$$\xi_{t,j}^h \leftarrow Emb_i(x_{t,j}^h)$$

**Step 2:** Take a weighted average of the embedding vectors:

$$\tilde{\xi}_t^h = [\xi_{t,1}^h \quad \xi_{t,2}^h \quad \dots \quad \xi_{t,m_2}^h] Softmax\left(W_2 [\xi_{t,1}^{hT} \quad \xi_{t,2}^{hT} \quad \dots \quad \xi_{t,m_2}^{hT}]^T + \mathbf{b}_2\right)$$

**Step 3:** Convert the temporal label sequence  $TimeSeq_t$  into a discrete feature vector, and add it to the representation vector:

$$\xi_t^h = \tilde{\xi}_t^h + Emb(TimeSeq_t)$$

## 375 A.2 LABEL LEARNING NETWORK

376 The Label Learning is designed to obtain the multi step predictions of the sequence label values.

## 377 A.2.1 SEQ2SEQ MODULE

378 As shown in Figure 6, the Seq2seq module is composed of an Encoder module and a Decoder  
 379 module. Each module corresponds to a LSTM cell and its subsequent network structure. Similar  
 380 to the Representation Learning, the Seq2seq module also needs to deal with the static features and  
 381 sequence features. In addition, we also need to add the position encoding at the input side of the  
 382 Decoder. The encoded position vector is

$$\xi_{t+j}^f = Emb(TimeSeq_{t+j}),$$

383 where  $TimeSeq_{t+j} = (t - p, t - p + 1, \dots, t, \dots, t + j)$ .

384 Suppose the representation vectors of the static and sequence features at time  $t$  from the Representa-  
 385 tion Learning module is  $\xi^s$  and  $\xi_t^h$ . We use  $\xi^s$  to initialize the hidden state  $\mathbf{h}^{en}$  and the unit state  
 386  $\mathbf{c}^{en}$  of the Encoder, where

$$\begin{aligned} \mathbf{h}_{t-p-1}^{en} &= GLU(\xi^s) = (W_3 \xi^s + b_3) \times \sigma(W_4 \xi^s + \mathbf{b}_4), \\ \mathbf{c}_{t-p-1}^{en} &= GLU(\xi^s) = (W_5 \xi^s + b_5) \times \sigma(W_6 \xi^s + \mathbf{b}_6), \end{aligned}$$

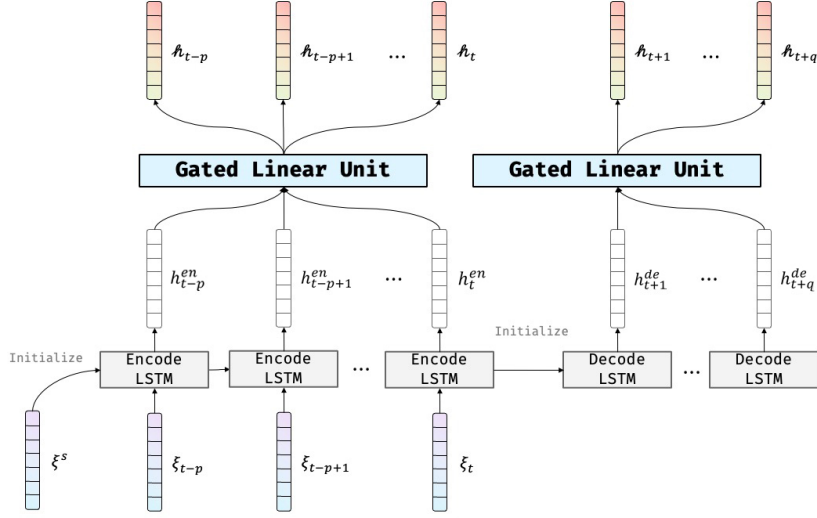


Figure 6: The operational structure of the Seq2seq module

387 where  $\sigma(x) = \frac{1}{1+\exp(-x)}$  is the sigmoid function. Then the Encoder model can calculate the hidden  
 388 state of the current time  $\mathbf{h}_t^{en}$  and the cell state  $\mathbf{c}_t^{en}$ .  $\mathbf{h}_t^{en}$  is the output of the Encoder at time  $t$ . This  
 389 process is different from the common Auto Encoder algorithm and the Seq2seq algorithm, where  
 390 the common Seq2seq algorithm only cares about the last hidden state of the Encoder and the output  
 391 of the Decoder. The main purpose of our Seq2seq module is to store the results of each step and  
 392 pass them to the subsequent structures for further feature extractions of the timing characteristics of  
 393 the data.

394 According to the model equation 1, the dimension of the input matrix of the Encoder is  $(p+1) \times d_k$ .  
 395 The GLU (Gated Linear Unit) unit is a shallow back propagation neural network, which is often used  
 396 for transition operations between the output and input of different structures. The computational  
 397 process of the Encoder is

$$\begin{aligned} \{\mathbf{h}_{t-i}^{en}, \mathbf{c}_{t-i}^{en}\} &= LSTM(\boldsymbol{\xi}_{t-i}^h, \mathbf{h}_{t-i-1}^{en}, \mathbf{c}_{t-i-1}^{en}), \\ \tilde{\mathbf{h}}_{t-i} &\leftarrow (W_7 \mathbf{h}_{t-i}^{en} + \mathbf{b}_7) \times \sigma(W_8 \mathbf{h}_{t-i}^{en} + \mathbf{b}_8), \\ \mathbf{h}_{t-i} &\leftarrow \tilde{\mathbf{h}}_{t-i} + W_{glu1} \tilde{\mathbf{h}}_{t-i}, \end{aligned}$$

398 where  $i = p, p-1, \dots, 0$ .

399 Similarly, on the Decoder side, the hidden state  $\mathbf{h}_t^{de}$  is the output of the Decoder at time  $t$ . The  
 400 Decoder uses the last hidden state and unit state of the Encoder to initialize the hidden state and unit  
 401 state. According to the model equation 1, the length of the input sequence at the Decoder side is  $q$ .  
 402 The calculation process of the Decoder is

$$\begin{aligned} \mathbf{h}_t^{de} &= \mathbf{h}_t^{en}, \quad \mathbf{c}_t^{de} = \mathbf{c}_t^{en} \\ \{\mathbf{h}_{t+j}^{de}, \mathbf{c}_{t+j}^{de}\} &= LSTM(\boldsymbol{\xi}_{t+j}^f, \mathbf{h}_{t+j-1}^{de}, \mathbf{c}_{t+j-1}^{de}), \\ \tilde{\mathbf{h}}_{t+j} &\leftarrow (W_9 \mathbf{h}_{t+j}^{de} + \mathbf{b}_9) \times \sigma(W_{10} \mathbf{h}_{t+j}^{de} + \mathbf{b}_{10}) \\ \mathbf{h}_{t+j} &\leftarrow \tilde{\mathbf{h}}_{t+j} + W_{glu2} \tilde{\mathbf{h}}_{t+j} \end{aligned}$$

403 where  $j = 1, 2, \dots, q$ .

404 Algorithm 4 is the computational workflow of our Seq2seq module. Our Seq2seq module adopts  
 405 some common structures, such as GLU and AddNorm, to improve the model training performance.  
 406 The AddNorm layer includes a Skip Connection structure (He et al., 2016) and a Batch Normal-  
 407 ization structure (Ioffe and Szegedy, 2015). The Skip Connection structure is used to reduce the  
 408 non-convexity of the network to protect the model from the gradient problem caused by the deep  
 409 learning structure. The Batch Normalization improves the training speed by normalizing each fea-  
 410 ture of samples within each batch. It is suitable for processing the structural data which has relatively  
 411 strong features independencies. The output of the Seq2seq module will be passed to the Attention  
 412 unit for the further analysis.

---

**Algorithm 4: Seq2seq Module**


---

**Input:**  $\{\xi^s, \xi_{t-p}^h, \xi_{t-p+1}^h, \dots, \xi_t^h\}$  the embedding vectors from the Representation Learning

**Output:**  $\{\mathbf{h}_{t-p}, \mathbf{h}_{t-p+1}, \dots, \mathbf{h}_t, \dots, \mathbf{h}_{t+q}\}$

**The Encoder stage:** Initialize the hidden state and unit state of the Encoder LSTM  $\mathbf{h}_{t-p-1}^{en}$   
 and  $\mathbf{c}_{t-p-1}^{en}$ ,

$$\mathbf{h}_{t-p-1}^{en} \leftarrow (W_3 \xi^s + b_3) \times \sigma(W_4 \xi^s + \mathbf{b}_4)$$

$$\mathbf{c}_{t-p-1}^{en} \leftarrow (W_5 \xi^s + b_5) \times \sigma(W_6 \xi^s + \mathbf{b}_6)$$

**For**  $i = p : 0$ , compute:

$$\{\mathbf{h}_{t-i}^{en}, \mathbf{c}_{t-i}^{en}\} \leftarrow LSTM(\xi_{t-i}^h, \mathbf{h}_{t-i-1}^{en}, \mathbf{c}_{t-i-1}^{en})$$

$$\tilde{\mathbf{h}}_{t-i} \leftarrow (W_7 \mathbf{h}_{t-i}^{en} + \mathbf{b}_7) \times \sigma(W_8 \mathbf{h}_{t-i}^{en} + \mathbf{b}_8)$$

$$\mathbf{h}_{t-i} \leftarrow \tilde{\mathbf{h}}_{t-i} + W_{glu1} \tilde{\mathbf{h}}_{t-i}$$

**The Decoder stage:** Initialize the hidden state and unit state of the Decoder LSTM unit  $\mathbf{h}_t^{de}$   
 and  $\mathbf{c}_t^{de}$ ,

$$\mathbf{h}_t^{de} \leftarrow \mathbf{h}_t^{en}, \mathbf{c}_t^{de} \leftarrow \mathbf{c}_t^{en}$$

$$\xi_{t+j}^f = Emb(TimeSeq_{t+j})$$

**For**  $j = 1 : q$ , compute:

$$\{\mathbf{h}_{t+j}^{de}, \mathbf{c}_{t+j}^{de}\} \leftarrow LSTM(\xi_{t+j}^f, \mathbf{h}_{t+j-1}^{de}, \mathbf{c}_{t+j-1}^{de})$$

$$\tilde{\mathbf{h}}_{t+j} \leftarrow (W_9 \mathbf{h}_{t+j}^{de} + \mathbf{b}_9) \times \sigma(W_{10} \mathbf{h}_{t+j}^{de} + \mathbf{b}_{10})$$

$$\mathbf{h}_{t+j} \leftarrow \tilde{\mathbf{h}}_{t+j} + W_{glu2} \tilde{\mathbf{h}}_{t+j}$$


---

### 413 A.2.2 ATTENTION MODULE

414 The Attention module analyze the sequence information processed by the Seq2seq module. Figure 7  
 415 is the operational structure diagram of the Attention module. First, the Attention module maps the  
 416 output of the Encoder-Decoder  $\{\mathbf{h}_{t+i} \mid -p \leq i \leq q\}$  and the static feature vector  $\xi^s$  into the input  
 417 matrix of the Attention module, where the input matrix  $M$  is defined as

$$M = [\mathbf{m}_{t-p} \ \mathbf{m}_{t-p+1} \ \dots \ \mathbf{m}_{t+q}]^T,$$

$$\mathbf{m}_{t+i} = W_{11}^{(1)} \mathbf{h}_{t+i} + W_{11}^{(2)} \xi^s + \mathbf{b}_{11},$$

418 where  $-p \leq i \leq q$ . Since the proposed algorithm is composed of the self-attention unit,  $M$  can be  
 419 used as the  $K$  and  $V$  matrices in the Attention layer, which means we can set  $K = V = M$ . For the  
 420  $q$ -step prediction of time series problem studied in this article, we can only use the corresponding  
 421 input on the Encoder side as the source of  $Q$  matrix, and then do the following calculation:

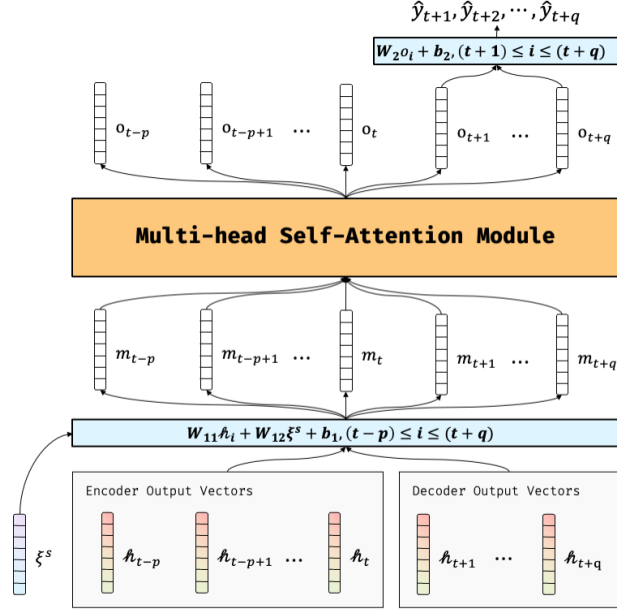


Figure 7: The operational structure of the Attention module

422 For the  $q$ -step predictions of the time series problem studied in this article, we only need to calculate  
 423 the latent vectors from the moment  $t + 1$  to  $t + q$  through the Attention layer. Therefore  $Q$  is  
 424 constructed in the following way,

$$Q = [\mathbf{m}_{t+1} \ \mathbf{m}_{t+2} \ \dots \ \mathbf{m}_{t+q}]^T.$$

425 In this case, the result of the Attention output is the hidden vectors from time  $t + 1$  to  $t + q$ . We can  
 426 also replace  $Q$  by  $M$ , then the output of the Attention layer will be the latent vectors from time  $t - p$   
 427 to  $t + q$ , which will be a waste of the computing resources since we only need the hidden vectors  
 428 from time  $t + 1$  to  $t + q$ . The final output of the Attention layer is the weighted average of the results  
 429 of each self-attention head.

$$O_h \leftarrow \text{Softmax}\left(\frac{W_h^q Q (W_h^k K)^T}{\sqrt{d_k}}\right) W_h^v V,$$

$$O \leftarrow \sum_{h=1}^H O_h, \quad O \in \mathbb{R}^{q \times d_k},$$

430 where  $H$  is the number of heads of the Self-attention module,  $W_h^q, W_h^k, W_h^v$  are the parameter  
 431 matrices of each self-attention head. The output of the self-attention mechanism is also processed  
 432 by the GLU and AddNorm units. These two operations do not change the dimension of the vectors.  
 433 The output of the Attention layer is then processed by the Position-Wise Feed Forward Network  
 434 Module to obtain the multi-step prediction vector  $\mathbf{y}_t^{(q)}$  for the next  $q$  moments in the future. The  
 435 computational process is

$$\begin{aligned} \tilde{O}_1 &\leftarrow O + (OW_{12} + \mathbf{b}_{12}) \times \sigma(OW_{13} + \mathbf{b}_{13}), \\ \tilde{O}_2 &\leftarrow [\mathbf{h}_{t+1} \ \mathbf{h}_{t+2} \ \dots \ \mathbf{h}_{t+q}]^T + (\tilde{O}_1 W_{14} + \mathbf{b}_{14}) \times \sigma(\tilde{O}_1 W_{15} + \mathbf{b}_{15}), \\ \mathbf{y}_t^{(q)} &= W_{out} \tilde{O}_2 + \mathbf{b}_{out}. \end{aligned}$$

436 This network structure comes from the Transformer. It maps the high dimensional vector corre-  
 437 sponding to each time step into a scalar through a shallow network. The scalar is the predicted label  
 438 value. Algorithm 5 is the computational workflow of Attention module.



---

**Algorithm 5:** Attention Module

---

**Input:** The output vector of the Seq2seq module  $(\mathbf{h}_{t-p}, \mathbf{h}_{t-p+1}, \dots, \mathbf{h}_t, \dots, \mathbf{h}_{t+q})$ **Output:**  $\mathbf{y}_t^{(q)}$  the vector of the  $q$ -step predictions from moment  $t + 1$  to  $t + q$ **For**  $i = -p: q$ , compute:

$$\mathbf{m}_{t+i} \leftarrow W_{11}^{(1)} \mathbf{h}_{t+i} + W_{11}^{(2)} \boldsymbol{\xi}^s + \mathbf{b}_{11}$$

the parameter matrix of the self-attention mechanism:

$$Q \leftarrow [\mathbf{m}_{t+1} \ \mathbf{m}_{t+2} \ \dots \ \mathbf{m}_{t+q}]^T$$

$$K = V = M \leftarrow [\mathbf{m}_{t-p} \ \mathbf{m}_{t-p+1} \ \dots \ \mathbf{m}_{t+q}]^T.$$

**For**  $h=1: H$ , compute:

$$O_h \leftarrow \text{Softmax}\left(\frac{W_h^q Q (W_h^k K)^T}{\sqrt{d_k}}\right) W_h^v V$$

Aggregate and compute the output of each attention mechanism head:

$$O \leftarrow \sum_{h=1}^H O_h$$

$$\tilde{O}_1 \leftarrow O + (O W_{12} + \mathbf{b}_{12}) \times \sigma(O W_{13} + \mathbf{b}_{13})$$

$$\tilde{O}_2 \leftarrow [\mathbf{h}_{t+1} \ \mathbf{h}_{t-2} \ \dots \ \mathbf{h}_{t+q}]^T + (\tilde{O}_1 W_{14} + \mathbf{b}_{14}) \times \sigma(\tilde{O}_1 W_{15} + \mathbf{b}_{15})$$

Predict the label value for the next  $q$  steps in the sequence:

$$\mathbf{y}_t^{(q)} \leftarrow W_{out} \tilde{O}_2 + \mathbf{b}_{out}$$


---