
End to End Time Series Analysis

Anonymous Author(s)

Affiliation

Address

email

Abstract

Time series plays an important role in practical applications, including stock predictions, traffic scheduling, clinical diagnosis. Such applications face two challenges: on one hand, most relationships are nonlinear, on the other hand, sampled data could always be noisy and incomplete. In the last few years, the newly developed deep learning based approaches started appearing in the community and have shown superior performances than the traditional models. This work gives a review of the most popular deep learning models and other algorithms of end-to-end learning for time series problems. While the neural networks have shown promising results on static data, such as natural language processing and computer vision, applying them with respect to time is still a challenging work. An overall review is provided to discuss the details of the models, learning features and tasks.

1 Introduction

With widespread adoption of electronic records, there is an increasing availability of large amounts of time series data and the need of performing accurate forecasting of future behavior in several scientific and applied domains. Time series usually contain temporal dependencies which may cause two otherwise identical samples to have different behavior. This characteristic demands the definition of robust and efficient techniques able to infer from observations the stochastic dependency between past and future. The forecasting domain used to be dominated by linear statistical methods such as ARIMA, which requires complete data and has to impute missing values before prediction[2][1]. More recently, with the possibility to model historical data as a nonlinear function and approximate complex latent relationships, machine learning models have established themselves as serious contenders to classical statistical models[5][4][3]. Here we give a review of those end-to-end (E2E) models for time series problems, where E2E refers to training a possibly complex learning system represented by a single model (specifically a Deep Neural Network) that represents the complete target system, bypassing the intermediate layers usually present in traditional pipeline designs.

In this work, we mainly investigate various end-to-end time series tasks: traditional time series forecasting, i.e., predicting the trend of a stock, time series classification (TSC), where we address to build a classifier to label the time series samples, and time series clustering. Formally, in the area of time series problems, given a sequence $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\} = X$ of observed states, where each element $\mathbf{x}^{(t)} \in \mathbb{R}^m$ is a vector describing the physical status at time step t such that $\{x_1^{(t)}, \dots, x_m^{(t)}\} = \mathbf{x}^{(t)}$, the E2E problem aims to predict the value of or label $\mathbf{x}^{(n+\Delta)}$ denoted by $\hat{\mathbf{y}}$ for some prediction horizon Δ [6][7], which is usually a fixed value in most scenarios.

The rest of this paper is structured as follows: Section 2 introduces basic types of neural network modules that are commonly used for time series tasks and for complicated model build-ups. Section 3, 4, 5 describe some tasks using end-to-end deep learning to perform modeling, classification and clustering. At the end, Section 6 concludes this work.

2 Supervised Learning Setting

In end-to-end time series analysis, supervised learning consists in modeling, on the basis of a finite set of observations, the relation between a set of input variables and one or more output variables, which are considered somewhat dependent on the inputs. To model an input/output mapping, i.e. to achieve E2E learning, the general approaches rely on the availability of a collection of data points typically referred to as training set that we denoted early by X , as well as the output value Y .

3 Neural Networks

For end-to-end time series, there are countless variations of deep learning models and it would not be possible to cover all the proposed methods. Hence, for simplicity, we focus on three basic types of neural networks known as *Multilayer Perceptrons (MLP)*, *Convolutional Neural Networks (CNN)* and *Recurrent Neural Networks (RNN)*, as well as some additional ones developed from those models. The reason for generally selecting the following models is that they enjoy the popularity along with the lucidity in the community.

3.1 Multilayer Perceptrons

MLPs are one of the first proposed architectures in artificial neural networks (ANN) and are considered to be the simplest type of neural network architectures. Though there are a significant number of studies on MLPs for time series problems, MLPs consist of three mainly components: input layer, hidden layer and output layer.

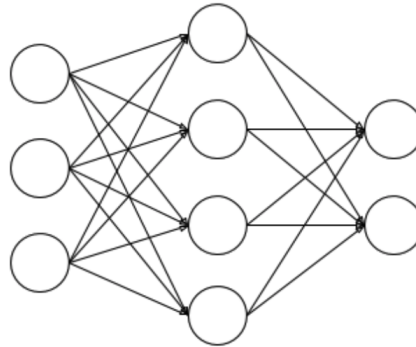


Figure 1: The basic structure of multilayer perceptrons. Image via [UCB EECS189](#).

Each node (neuron) in the hidden layers and in the output layers computes a weighted sum of its inputs \mathbf{x} with the weight \mathbf{w} and the bias term b . To approximate nonlinear relationships, on top of that, there is a nonlinear activation function σ for the nodes that is usually layer-wise shared. Concretely, the computation on the node can be expressed as

$$\sigma(\mathbf{w}^\top \mathbf{x} + b).$$

Let L be the number of layers in an MLP, then the entire model computes the function

$$\sigma_L(\mathbf{W}_L \sigma_{L-1}(\cdots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) \cdots) + \mathbf{b}_L),$$

where \mathbf{W} and \mathbf{b} are the matrix of weights and the vector of bias in the layer l respectively.

The important issues for MLPs are the hyperparameters, the activation functions and the training method, which have been mentioned in the previous lectures of the course and will not be repeated here. For time series forecasting, previous studies MLPs can be powerful and universal due to its good generalization ability. However, with a comparatively simpler structure, in most cases MLP perform worse than other complex networks, e.g. LSTM [8][9].

3.2 Convolutional Neural Networks

CNN is a type of neural networks that consists some specific layers (convolutional and pooling layers) for convolutional operations. CNN enjoys all the benefits that MLP has while not requiring to learn from lag observations which MLP does, while significantly reducing the number of weights. Though having been used in various fields, like object segmentation, style conversion, automatic coloring, what CNN can really do is just function as a feature extractor.

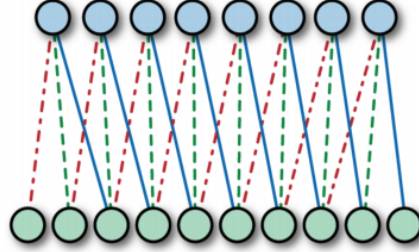


Figure 2: The basic structure of convolutional layers. Image via [UCB EECS189](#).

A convolutional layer takes a $W \times H \times D$ dimensional input I and convolves it with a $w \times h \times D$ dimensional filter G . Mathematically, the convolution operator is defined as

$$(\mathbf{I} * \mathbf{G})[x, y] = \sum_{a=0}^{w-1} \sum_{b=0}^{h-1} \sum_{c \in \{1 \dots D\}} I_c[x + a, y + b] \cdot G_c[a, b].$$

From the operation of convolution, it can be seen that any unit in the output image is only related to a part of the input image. In traditional neural networks, since they are all fully connected, any unit of output must be affected by all inputs. In this way, the recognition effect of the potential features will be greatly reduced. Specifically speaking, each area has its own unique characteristics, which we wish not to be affected by other areas.

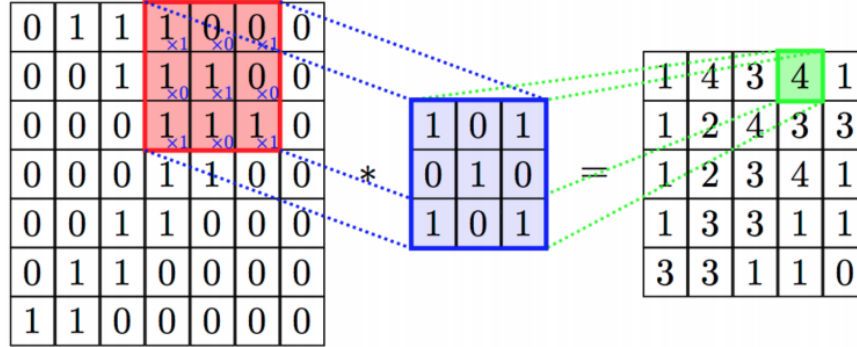


Figure 3: Convolving a filter with an image, where $W = H = 7$, $D = 1$ and $w = h = 3$. Image via [UCB EECS189](#).

3.3 Recurrent Neural Networks

The most popular model that have been used for modeling sequential data is the Recurrent Neural Network [10][11]. In RNN, the activations from each time are stored in the internal state of the model in order to provide the memory at that moment.

Had there not been the cycle arrow in the left part of fig. 4, the wrapped RNN would be an MLP, which is made up by the input, the hidden component and the output. Let U , V and W be the weight matrices between a input state and the hidden state (x_t and h_t), between a hidden state and the output state at the same time step (h_t and y_t), and between hidden states (h_{t-1} and h_t) respectively. With

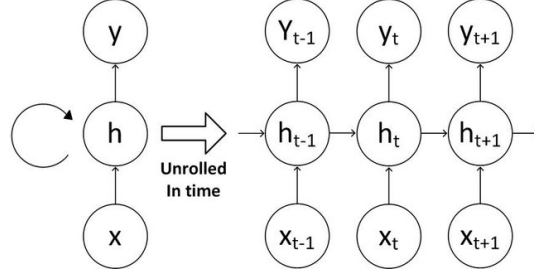


Figure 4: A conceptual visualization of the recurrent nature of an RNN, where x is the input data for the input layer, h represents the hidden layers that carry internal state (memory) through the time, and y is the model output. Image by Savvas Varsamopoulos.

b and σ remain to denote the bias term between hidden layers and the activation function severally, additionally letting c be the bias term between the hidden layer and the output layer, we have the equations in the computational graph as:

$$\begin{aligned} h_t &= \sigma(Ux_t + Wh_{t-1} + b) \\ y_t &= \sigma(Vh_t + c) \end{aligned}$$

Note that σ may differ regarding disparate targets of the network structures.

Though theoretically RNN works perfectly as a sequential data approximator, the major challenge with a typical generic RNN is that these networks remember only a few earlier steps in the sequence and thus are not suitable when a longer sequence is involved. Specifically speaking, during the training the gradient of some weights starts to become too small or too large to compute, known as *gradient vanishing/gradient exploding*. A type of neural network architecture that solves the problem is Long Short-Term Memory.

3.4 Long Short-Term Memory

LSTM is a special kind of Recurrent Neural Network (RNN) with the capability of remembering the values from the past for the purpose of future use [12]. A key reason for the success of LSTM in time series problems is its ability to filter and keep memorization from the earlier stage by the recurrent connected units (*cells*) and the gates along with two horizontal lines incorporated, which carry the memory.

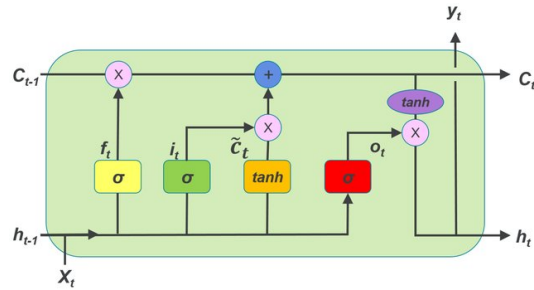


Figure 5: A cell of Long Short-Term Memory Neural Network, where the colored components are different components of the network. Image by Ismail et al.

Regarding a typical LSTM, there are three kinds of gates: forget, input and output. Based on fig. 5, we will describe its internal elements on an individual basis. which are surrounding and connecting by various of lines. First is the *forget gate*, the leftmost part of the cell, i.e. the yellow rectangle and the lines through it, which absorb the hidden state from the last cell $h^{(t-1)}$ as well as the newly input $X^{(t)}$. The sigmoid activate function output $f^{(t)} \in [0, 1]$, representing the extent to which the past memory remains in the current cell:

$$f^{(t)} = \sigma \left(W_f h^{(t-1)} + U_f x^{(t)} + b_f \right),$$

110 where U , W , b have the same meaning as the ones in RNN.

111 *Input gate*, built up by the green and the orange activation functions and the lines through them,
112 controls the extent to which $X^{(t)}$ would flow into the cell, which is expressed by two parts:

$$i^{(t)} = \sigma \left(W_i h^{(t-1)} + U_i x^{(t)} + b_i \right),$$

$$a^{(t)} = \tanh \left(W_a h^{(t-1)} + U_a x^{(t)} + b_a \right).$$

113 Before moving to the *Output gate*, we have to firstly take care of the current cell state $C^{(t)}$ and the
114 effects of the former two gates:

$$C^{(t)} = C^{(t-1)} \odot f^{(t)} + i^{(t)} \odot a^{(t)}.$$

115 Here \odot denotes the Hadamard product.

116 Then with the new born cell state, the rightmost part of the cell, i.e. the *output gate*, will generate the
117 hidden state $h^{(t)}$ and the prediction $\hat{y}^{(t)}$ as follows:

$$o^{(t)} = \sigma \left(W_o h^{(t-1)} + U_o x^{(t)} + b_o \right),$$

$$h^{(t)} = o^{(t)} \odot \tanh \left(C^{(t)} \right),$$

$$\hat{y}^{(t)} = \sigma \left(V h^{(t)} + c \right).$$

118 4 Time Series Forecasting

119 Almost all the neural networks mentioned above are suitable for performing time series forecasting
120 tasks. Benefited from the availability of universal approximation, neural networks have almost
121 dominated this field. Several different deep learning approaches can be found in the literature for
122 performing forecasting tasks [13][14][15].

123 LSTM, by its nature utilizes the temporal characteristics of any time series signal, hence forecasting
124 end-to-end time series is a well-studied and successful implementation of LSTM. Meanwhile, series
125 combination models, especially linear/nonlinear hybrid models (e.g. ARIMA/LSTM), are not
126 uncommon in the literature [16][17].

127 When it come to the metrics of the experiments, besides the MAE and MSE that have been mentioned
128 in the previous lectures, here we introduce a few more common metrics for time series analysis.

129 Root mean square error (RMSE) is the standard deviation of the prediction errors, can be expressed as

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}.$$

130 Mean absolute percentage error gives a good idea of the relative error, whose formula is given by

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \cdot 100\%.$$

131 This approach, however, is facing the problem when the sequences can have small denominators that
132 leads to a zero or nearly zero division. Since MAPE puts a heavy penalty on negative errors, as a
133 consequence, following this metric the system may prefer the model with a lower prediction. This
134 drawback can be avoided by taking the logarithm of the accuracy ratio.

135 Last, we have R-squared

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}},$$

136 where SS_{res} is the is the sum of squared residuals from the predicted values and SS_{tot} is the sum of
137 squared deviations of the dependent variable from the sample mean. As the fraction of the total sum
138 of squares that is explained by the model, this mechanic shows whether the model fits the observed
139 samples and how good it fits, where a high R^2 stands for a high correlation.

5 Time Series Classification

Time series classification (TSC), the problem to label the time series samples, has become a popular task of time series [18]. Given the need to accurately classify time series data, researchers have proposed hundreds of methods to solve this task, such as a nearest neighbor classifier with Dynamic Time Wrapping (DTW) [19][20], which has been proved to be a promising baseline.

Regarding end-to-end learning as the topic of this notes, here we introduce some CNN based methods from related studies [21][22][23]. After revolutionizing the performances of several tasks in computer vision like object classifications, face verifications and etc., CNN has been introduced to this topic to address TSC problems because of its success in representation learning. One of the recent works inspired by computer vision encodes the time series data as images and then apply CNN as a classifier [24], as CNN can treat the raw input data as a 1-D image then read it and store it as important element. Compared with k -NN with DTW, such a deep learning framework can learn a hierarchical feature representation from raw data automatically, matching the concept of end-to-end learning perfectly [25].

Here we talk about a work combining LSTM, convolutional network with attention mechanism: ALSTM-FCN [26].

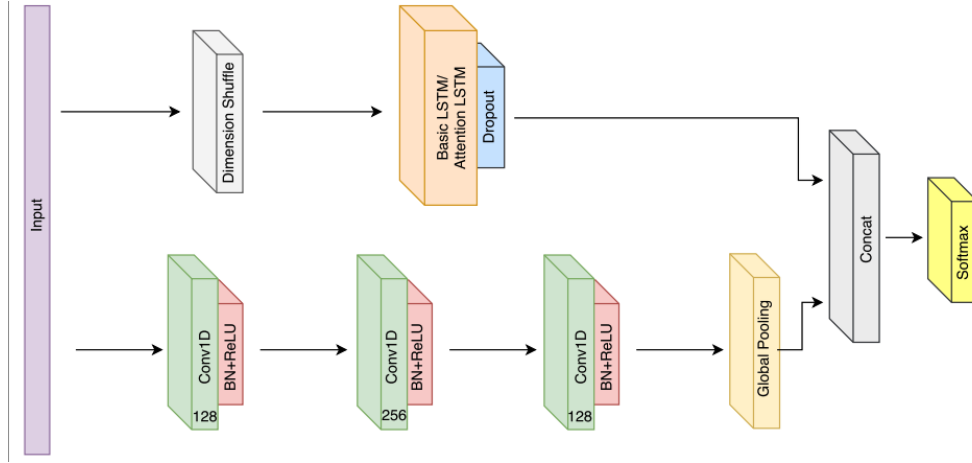


Figure 6: The LSTM-FCN architecture. LSTM cells can be replaced by Attention LSTM cells to construct the ALSTM-FCN architecture. Image by Karim et al.

The attention mechanism is a technique often used in neural translation of text, where a context vector C is conditioned on the target sequence y . The context vector c_i depends on a sequence of annotations $\{h^{(1)}, \dots, h^{(n)}\}$ to which an encoder maps the input sequence. Each annotation h_i contains information about the whole input sequence with a strong focus on the parts surrounding the i -th word of the input sequence. The context vector c_i is then computed as a weighted sum of these annotations h_i :

$$c_i = \sum_{j=1}^n \alpha_{ij} h_j$$

The weight α_{ij} of each annotation h_j is computed by:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^n \exp(e_{ik})}$$

where $e_{ij} = a(s_{i-1}, h_j)$ is an alignment model, which scores how well the input around position j and the output at position i match. Recent works have argued that attention mechanisms, without any recurrence, can be effective in end-to-end modeling tasks [27][28].

The fully convolutional block consists of three stacked temporal convolutional blocks with filter sizes of 128, 256, and 128 respectively. Each convolutional block is identical to the convolution block in

the CNN architecture. Each block consists of a temporal convolutional layer, which is accompanied by batch normalization and followed by a ReLU activation function. Finally, global average pooling is applied after the final convolution block.

Simultaneously, the time series input is conveyed into a dimension shuffle layer. The transformed time series from the dimension shuffle is then passed into the LSTM block. The LSTM block, comprising of either a general LSTM layer or an Attention LSTM layer, is followed by a dropout. The output of the global pooling layer and the LSTM block is concatenated and passed onto a softmax classification layer.

6 Anomaly Detection

An anomaly, or an outlier, is a sample that performing significantly different from the rest in the set, or in the sequence for temporal samples. Anomaly detection refers to retrieve, detect even predict those abnormal issues in time series sequences. Reliable temporal uncertainty estimation is a critical industrial demand, e.g. business and properties monitoring for Yahoo and Microsoft [29][30]. Due to the complexity and the comprehensiveness of the designs, a business supporting anomaly detection system always showed a far more complicated framework. In this section we will briefly introduce some representative systems.

6.1 Extensible GenericAnomaly Detection System (EGADS)

Proposed by a team from Yahoo, EGADS is the first comprehensive system for anomaly detection that is flexible, accurate, scalable and extensible. Though it is not an end-to-end learning model, we still put it here as it properly define the basic concepts for anomaly detection systems and is open source¹.

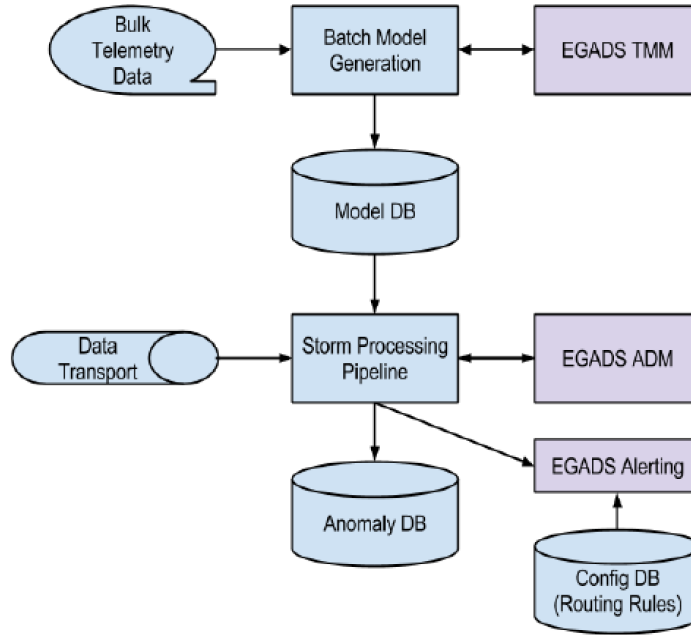


Figure 7: The architecture for EGADS-YMS. Image by Laptev et al.

Some recent works based on deep learning have compared themselves against EGADS and have made remarkable improvements on both real and synthetic data [31].

DeepAD, an anomaly detection model that leverages a plethora of basic models, has three main phases as illustrated in fig. 8: time series forecasting, merge predictions, and anomaly detector. At the first phase both a statistical model (ARIMA) and a nonlinear model (LSTM) will be trained, the

¹Source code available at <https://github.com/yahoo/egads>.

194 results of which will then be applied and combined to mix the advantages of those basic models.
 195 Then the anomaly detector ensued, computing a dynamic threshold at each time step on the past
 196 scaled squared error. This is a generic anomaly detection framework that does not utilize the prior
 197 knowledge neither for training nor for decision making.

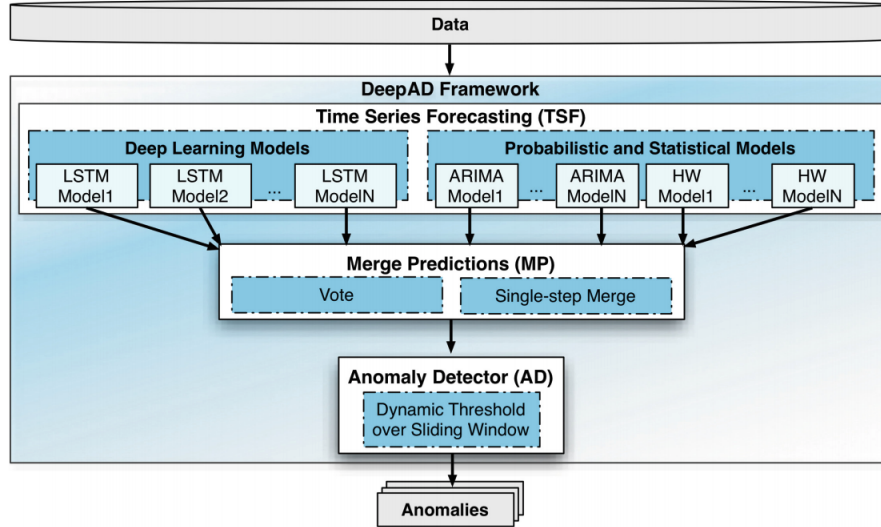


Figure 8: The framework overview of DeepAD. Image by Buda et al.

198 Meanwhile, there are some studies focusing on simpler scenarios, e.g. holiday prediction for Uber
 199 trips [32][33]. In the following work, the Uber team build their model with an encoder-decoder
 200 framework for inherent temporal pattern capturing and a prediction network that takes both the
 201 encoder-decoder output and potential external features to guide the final prediction.

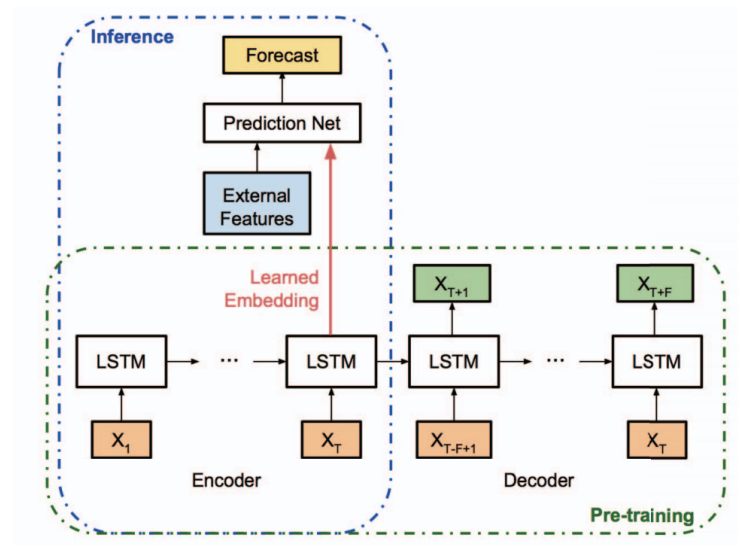


Figure 9: The neural network architecture proposed by Uber team. Image by Zhu and Laptev.

7 Prophet

Besides traditional models like ARIMA and those fancy neural networks, there is also a sophisticated yet multipurpose tool provided by Facebook for time series forecasting, Prophet[34]. Based on time series decomposition and machine learning fitting, Prophet can not only handle the case of some outliers in the time series problems, but also the case of some missing values, and predict the future trend of the time series almost automatically. This tool that provides programming interfaces in both R and Python enables statisticians and analysts to use their expertise for business or research purposes. In this section, we will briefly introduce the methodology of Prophet and leave some practice problems in the homework.

7.1 Algorithm

In the field of traditional time series analysis, there is a common analysis method called time series decomposition, where let y_t be the time series. Then we can write

$$y_t = S_t + T_t + R_t$$

where S_t is the seasonal component, T_t is the trend component, and R_t is the remainder component, all with respect to period t . Similarly, a multiplicative decomposition would be written as

$$y_t = S_t \times T_t \times R_t,$$

which is equivalent to

$$\ln y_t = \ln S_t + \ln T_t + \ln R_t.$$

Generally speaking, in real life and business tasks, in addition to seasonality, trend, and remainders, there are usually holiday effects, which in Prophet are expressed as

$$y(t) = g(t) + s(t) + h(t) + \varepsilon_t,$$

where $g(t)$ is the piecewise linear or logistic growth curve for modelling non-periodic changes in time series, $s(t)$ is the periodic changes (e.g. weekly/yearly seasonality), $h(t)$ is the effects of holidays (user provided) with irregular schedules, and $\varepsilon(t)$ is the error term accounts for any unusual changes not accommodated by the model.

7.2 Application

As a competitive time series analysis tool, Prophet has been applied for real world business predictions, where the data are highly non stationary and irregular. Here we will focus on a work in retail industry, which presents a framework capable of accurately forecasting future sales in the retail industry and classifying the product portfolio according to the expected level of forecasting reliability [35]. The proposed framework, that would be of great use for any company operating in the retail industry, is based on Facebook's Prophet algorithm and backtesting strategy. Real-world sales forecasting benchmark data obtained experimentally in a production environment in one of the biggest retail companies in Bosnia and Herzegovina is used to evaluate the framework and demonstrate its capabilities in a real-world use case scenario.

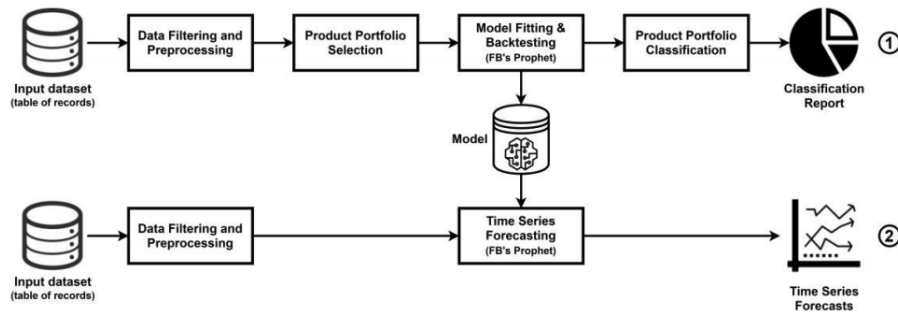


Figure 10: The proposed sales forecasting model. Image by Žunić et al.

In the work flow shown in fig. 10, after data preprocessing, Prophet is used for modelling the dynamics of sales for items in a product portfolio without using additional regressors, with the aim of generating monthly and quarterly sales forecasts. It is empirically concluded that at least 24 months of historical data is required for reliable estimation of trend and/or seasonal effects.

8 Conclusions

In addition to the above ones, there are also other interesting end-to-end works for time series tasks: an imputation network with residual short paths for simultaneous time series data imputation and prediction [36], Apply *Gated Restricted Boltzmann Machine* (GRBM) for video action recognition[37], and etc.

The dominance of RNN-based models for end-to-end time series analysis will probably not disappear anytime soon, mainly due to their easy adaptation to most temporal dependency problems. Meanwhile, some enhanced versions of the original LSTM or RNN models, generally integrated with hybrid learning systems started becoming more common. Here we grouped the studies and listed some symbolic models. However, it indicates that there are still latent patterns and opportunities waiting for our readers to explore because of the high uncertainty of time series.

Appendix A Datasets

Here we offer some datasets that are widely used by the community.

- UCR Time Series Classification Archive: A superset for 48 datasets (updated in 2018 fall) for time series classification and clustering. Available at https://www.cs.ucr.edu/~eamonn/time_series_data_2018/.
- Human Activity Recognition Using Smartphones Data Set: A dataset recording the activities while using smartphones, where the features are obtained by calculating variables from the time and frequency domain. Available at <https://archive.ics.uci.edu/ml/datasets/human+activity+recognition+using+smartphones#>.
- Yahoo Webscope Program: A reference library of interesting and scientifically useful datasets, used for anomaly detection. Available at <https://webscope.sandbox.yahoo.com/>.

References

- [1] Hillmer, Steven Craig, and George C. Tiao. "An ARIMA-model-based approach to seasonal adjustment." *Journal of the American Statistical Association* 77.377 (1982): 63-70.
- [2] De Gooijer, Jan G., and Rob J. Hyndman. "25 years of time series forecasting." *International journal of forecasting* 22.3 (2006): 443-473.
- [3] Siami-Namini, Sima, Neda Tavakoli, and Akbar Siami Namin. "A comparison of ARIMA and LSTM in forecasting time series." 2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA). IEEE, 2018.
- [4] Bontempi, Gianluca, Souhaib Ben Taieb, and Yann-Aël Le Borgne. "Machine learning strategies for time series forecasting." *European business intelligence summer school*. Springer, Berlin, Heidelberg, 2012.
- [5] Ahmed, Nesreen K., et al. "An empirical comparison of machine learning models for time series forecasting." *Econometric Reviews* 29.5-6 (2010): 594-621.
- [6] Lippi, Marco, Matteo Bertini, and Paolo Frasconi. "Short-term traffic flow forecasting: An experimental comparison of time-series analysis and supervised learning." *IEEE Transactions on Intelligent Transportation Systems* 14.2 (2013): 871-882.
- [7] Venkatraman, Arun, Martial Hebert, and J. Andrew Bagnell. "Improving multi-step prediction of learned time series models." *Twenty-Ninth AAAI Conference on Artificial Intelligence*. 2015.
- [8] Chen, Lin, et al. "Which artificial intelligence algorithm better predicts the Chinese stock market?." *IEEE Access* 6 (2018): 48625-48633.

[9] Cao, Jian, Zhi Li, and Jian Li. "Financial time series forecasting model based on CEEMDAN and LSTM." *Physica A: Statistical Mechanics and its Applications* 519 (2019): 127-139.

[10] Hüsken, Michael, and Peter Stagge. "Recurrent neural networks for time series classification." *Neurocomputing* 50 (2003): 223-235.

[11] Sagheer, Alaa, and Mostafa Kotb. "Time series forecasting of petroleum production using deep LSTM recurrent networks." *Neurocomputing* 323 (2019): 203-213.

[12] Gers, Felix A., Douglas Eck, and Jürgen Schmidhuber. "Applying LSTM to time series predictable through time-window approaches." *Neural Nets WIRN Vietri-01*. Springer, London, 2002. 193-200.

[13] Romeu, Pablo, et al. "Time-series forecasting of indoor temperature using pre-trained deep neural networks." *International conference on artificial neural networks*. Springer, Berlin, Heidelberg, 2013.

[14] Turner, Jeffrey T. *Time series analysis using deep feed forward neural networks*. University of Maryland, Baltimore County, 2014.

[15] Sezer, Omer Berat, Mehmet Ugur Gudelek, and Ahmet Murat Ozbayoglu. "Financial time series forecasting with deep learning: A systematic literature review: 2005–2019." *Applied Soft Computing* 90 (2020): 106181.

[16] Khashei, Mehdi, and Zahra Hajirahimi. "A comparative study of series arima/mlp hybrid models for stock price forecasting." *Communications in Statistics-Simulation and Computation* 48.9 (2019): 2625-2640.

[17] Chaâbane, Najeh. "A hybrid ARFIMA and neural network model for electricity price prediction." *International journal of electrical power & energy systems* 55 (2014): 187-194.

[18] Fawaz, Hassan Ismail, et al. "Deep learning for time series classification: a review." *Data Mining and Knowledge Discovery* 33.4 (2019): 917-963.

[19] Berndt, Donald J., and James Clifford. "Using dynamic time warping to find patterns in time series." *KDD workshop*. Vol. 10. No. 16. 1994.

[20] Bagnall, Anthony, et al. "The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances." *Data Mining and Knowledge Discovery* 31.3 (2017): 606-660.

[21] Zheng, Yi, et al. "Time series classification using multi-channels deep convolutional neural networks." *International Conference on Web-Age Information Management*. Springer, Cham, 2014.

[22] Cui, Zhicheng, Wenlin Chen, and Yixin Chen. "Multi-scale convolutional neural networks for time series classification." *arXiv preprint arXiv:1603.06995* (2016).

[23] Wang, Zhiguang, Weizhong Yan, and Tim Oates. "Time series classification from scratch with deep neural networks: A strong baseline." *2017 International joint conference on neural networks (IJCNN)*. IEEE, 2017.

[24] Wang, Zhiguang, and Tim Oates. "Encoding time series as images for visual inspection and classification using tiled convolutional neural networks." *Workshops at the twenty-ninth AAAI conference on artificial intelligence*. Vol. 1. 2015.

[25] Pelletier, Charlotte, Geoffrey I. Webb, and François Petitjean. "Temporal convolutional neural network for the classification of satellite image time series." *Remote Sensing* 11.5 (2019): 523.

[26] Karim, Fazle, et al. "LSTM fully convolutional networks for time series classification." *IEEE access* 6 (2017): 1662-1669.

[27] Qin, Yao, et al. "A dual-stage attention-based recurrent neural network for time series prediction." *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. 2017.

[28] Song, Huan, et al. "Attend and diagnose: Clinical time series analysis using attention models." *32nd AAAI Conference on Artificial Intelligence, AAAI 2018*. AAAI press, 2018.

[29] Laptev, Nikolay, Saeed Amizadeh, and Ian Flint. "Generic and scalable framework for automated time-series anomaly detection." *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 2015.

[30] Ren, Hansheng, et al. "Time-Series Anomaly Detection Service at Microsoft." *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2019.

[31] Buda, Teodora Sandra, Bora Caglayan, and Haytham Assem. "Deepad: A generic framework based on deep learning for time series anomaly detection." *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, Cham, 2018.

[32] Zhu, Lingxue, and Nikolay Laptev. "Deep and confident prediction for time series at uber." *2017 IEEE International Conference on Data Mining Workshops (ICDMW)*. IEEE, 2017.

- 330 [33] Laptev, Nikolay, et al. "Time-series extreme event forecasting with neural networks at uber." International
331 Conference on Machine Learning. Vol. 34. 2017.
- 332 [34] Taylor, Sean J., and Benjamin Letham. "Forecasting at scale." The American Statistician 72.1 (2018):
333 37-45.
- 334 [35] Žunić, Emir. "Application of Facebook's Prophet Algorithm for Successful Sales Forecasting Based on
335 Real-world Data." International Journal of Computer Science & Information Technology (IJCSIT) Vol 12 (2020).
- 336 [36] Shen, Lifeng, Qianli Ma, and Sen Li. "End-to-end time series imputation via residual short paths." Asian
337 Conference on Machine Learning. 2018.
- 338 [37] Memisevic, Roland, and Geoffrey Hinton. "Unsupervised learning of image transformations." 2007 IEEE
339 Conference on Computer Vision and Pattern Recognition. IEEE, 2007.