



# Transductive LSTM for time-series prediction: An application to weather forecasting

Zahra Karevan\*, Johan A.K. Suykens

ESAT-STADIUS, Kasteelpark Arenberg 10, 3001 Leuven, Belgium

## ARTICLE INFO

### Article history:

Received 9 May 2019

Received in revised form 1 November 2019

Accepted 30 December 2019

Available online 8 January 2020

### Keywords:

Transductive learning

Long short-term memory

Weather forecasting

## ABSTRACT

Long Short-Term Memory (LSTM) has shown significant performance on many real-world applications due to its ability to capture long-term dependencies. In this paper, we utilize LSTM to obtain a data-driven forecasting model for an application of weather forecasting. Moreover, we propose Transductive LSTM (T-LSTM) which exploits the local information in time-series prediction. In transductive learning, the samples in the test point vicinity are considered to have higher impact on fitting the model. In this study, a quadratic cost function is considered for the regression problem. Localizing the objective function is done by considering a weighted quadratic cost function at which point the samples in the neighborhood of the test point have larger weights. We investigate two weighting schemes based on the cosine similarity between the training samples and the test point. In order to assess the performance of the proposed method in different weather conditions, the experiments are conducted on two different time periods of a year. The results show that T-LSTM results in better performance in the prediction task.

© 2020 Elsevier Ltd. All rights reserved.

## 1. Introduction

Recurrent Neural Networks (RNNs) (Elman, 1990) are one of the most popular data-driven approaches used for time-series prediction. RNNs benefit from feedback loops which makes the output of RNN in one time step to be seen as an input in the subsequent one. Having information in a sequential form, the input of the RNN in each time step includes the information of the corresponding time step in the sequence and the previous information provided by a feedback loop. Despite the power of RNNs to understand the short-term dependencies, they have problems in capturing long-term ones due to the vanishing gradient problem (Bengio, Simard, & Frasconi, 1994). The vanishing gradient problem indicates that by using a gradient-based optimization in RNNs, the error gradients vanish in the earlier inputs after propagating through several steps. Consequently, the information of earlier steps are neglected in prediction or classification task. To avoid the vanishing gradient problem in RNNs, Hochreiter & Schmidhuber proposed Long Short-Term Memory (LSTM) (Hochreiter & Schmidhuber, 1997) which is able to capture long-term dependencies by utilizing a memory cell and gating method. LSTMs have been widely used and shown significant performance on different sequence learning problems such as video classification, machine translation and speech recognition (Graves, Mohamed, & Hinton, 2013;

Ng et al., 2015; Sutskever, Vinyals, & Le, 2014). Recently, there has been an increasing interest towards using LSTMs for time-series prediction such as air pollution forecasting, diagnosing based clinical data and traffic flow prediction (Freeman, Taylor, Gharabaghi, & Thé, 2018; Lipton, Kale, Elkan, & Wetzal, 2016; Tian & Pan, 2015).

Mostly, learning methods are based on inductive learning, where it is assumed that there is a function that maps the input to the output and the aim is to approximate the function in an optimal way. Learning based on induction results in a global model based on all training data points and the model can be used to predict or classify an unseen data point regardless of its position in the feature space. The idea of locality was first introduced by Vapnik in Vapnik (1992) where he discussed the trade off between the capacity of learning and the number of data points. Due to the possible differences in distribution of data in various regions of the feature space, it is meaningful to consider different learning capacities and properties for different areas in the feature space (Bottou & Vapnik, 1992).

In transductive learning (also known as local learning), the learning objective is more focused on obtaining good performance in the vicinity of the unseen data point. Note that to classify or predict an unseen test point, the label is assumed to be unknown while the feature vector is known. Based on the idea of transductive learning, the models are created such that their performance in the neighborhood of the test point in the feature space is optimized (Pang & Kasabov, 2004). One intuitive

\* Corresponding author.

E-mail address: [zahra.karevan@gmail.com](mailto:zahra.karevan@gmail.com) (Z. Karevan).

approach is a weighting scheme that considers weights for the data points in the training set based on their distance from the test point: the data points in the neighborhood of the test point obtain higher weights. As shown in [Bottou and Vapnik \(1992\)](#), these weights can be binary or real. If the weights are binary, one may utilize  $k$ -Nearest Neighbor ( $k$ -NN) approach to select part of the data points and create the model based on these samples. Nevertheless, in case of real values, each data point receives a weight based on its similarity with the test point. The closest sample gets the highest weight and the furthest one receives the lowest one. Such weighting scheme was used in the framework of Moving Least Squares by [Levin \(1998\)](#).

Transductive learning has been used in different types of problems. Do & Poulet proposed a parallel ensemble learning algorithm of random local Support Vector Machines which can help in dealing with some of the issues related to big data problems ([Do & Poulet, 2015](#)). Gilardi & Bengio utilized transductive learning for spatial data analysis ([Gilardi & Bengio, 2000](#)). Zhai et al. deployed transductive learning in a multi-view approach ([Zhai et al., 2018](#)). In several applications, transductive learning outperforms inductive learning in terms of accuracy ([Chen et al., 2018](#); [Elattar, Goulernas, & Wu, 2010](#); [Wu, Li, Meng, & Ngan, 2018](#)).

The weather system is known as a challenging complex system, and reliable weather forecasting has been the subject of many studies. State-of-the-art methods utilized Numerical Weather Prediction (NWP) which is a computationally intense method and demands thousands of processing units to have reliable forecasting models ([Bauer, Thorpe, & Brunet, 2015](#)). Recently, to avoid the shortcomings of NWP, data-driven approaches have become a major interest of the researchers to predict the future behavior of the weather system. In our previous studies, we investigated transductive learning as a data-driven approach for weather prediction in the frameworks of feature selection and function approximation ([Karevan, Feng, & Suykens, 2016](#); [Karevan & Suykens, 2018](#)). The results of these studies suggest that deploying transductive learning can improve the predictions performance.

In this study, we deploy LSTM for temperature prediction and propose Transductive LSTM (T-LSTM) by altering the cost function in the regression problem. In the proposed method, the impact of the training samples on the cost function is determined by their similarity to the test point. In this study, the weights of the training data points are defined based on their cosine similarity with the test point.

The remaining of this paper is as follows: first, in Section 2 we discuss LSTM architecture and how they can be used for regression problems. Then, in Section 4 we explain the proposed T-LSTM approach. Afterwards, the experimental results on an application of weather forecasting are shown in Section 5. Finally, the conclusion is provided.

## 2. Long short-term memory

The original idea of LSTM was initially proposed by [Hochreiter and Schmidhuber \(1997\)](#) and since then different approaches have been proposed to improve the performance of LSTM ([Cho et al., 2014](#); [Gers & Schmidhuber, 2000](#); [Graves & Schmidhuber, 2005](#); [Schmidhuber, Wierstra, Gagliolo, & Gomez, 2007](#)). In this study, we describe and deploy the popular architecture proposed by [Gers, Schraudolph, and Schmidhuber \(2002\)](#) and used in many other studies such as [Graves \(2013\)](#), [Graves et al. \(2013\)](#) and [Zaremba, Sutskever, and Vinyals \(2014\)](#).

LSTM utilizes a gating mechanism to control the information that has to be kept over time, the duration it has to be kept and the time that it can be read through the memory cell ([Bandara, Bergmeir, & Smyl, 2017](#)). In this paper we consider the LSTM

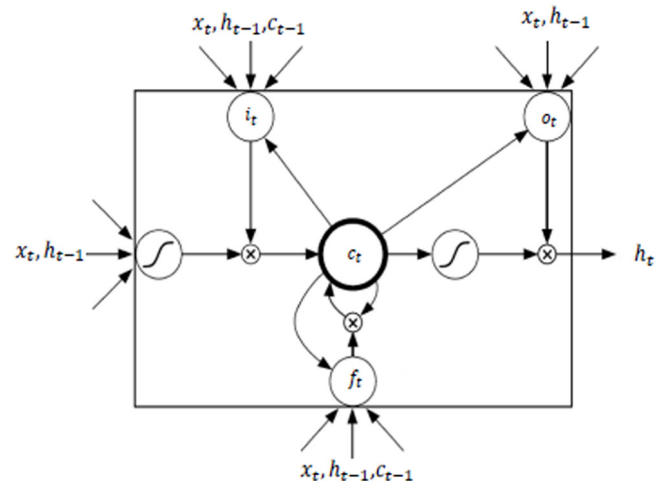


Fig. 1. LSTM cell based on [Graves \(2013\)](#).

cell as described in the paper of [Graves \(2013\)](#). To process the information, LSTM benefits from three gates. The architecture of an LSTM cell is depicted in [Fig. 1](#):

Assuming  $i_t, f_t, o_t, c_t$  and  $h_t$  to indicate the values of the input gate, forget gate, output gate, memory cell and hidden state at time  $t$  in the sequence respectively and  $x_t$  be the input of the system at time  $t$ , the architecture of the LSTM cell can be defined as follows ([Graves, 2013](#)):

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \quad (1)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \quad (2)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \quad (3)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \quad (4)$$

$$h_t = o_t \odot \tanh(c_t). \quad (5)$$

Note that  $\sigma(\cdot)$  represents the sigmoid as the activation function. Both logistic sigmoid and hyperbolic tangent are applied element-wise. The full weight matrices  $W_{xj}$  for  $j \in \{i, f, o, c\}$  are the weights for the input  $x_t$  in input, forget and output gates and the memory cell. The weight matrices  $W_{cj}$  for  $j \in \{i, f, o\}$  are diagonal matrices that are related to the connection of the cell memory to different gates. Note that the number of neurons for all gates is a predefined parameter and Eqs. (1) to (5) are applied for each neuron. Denoting  $n$  as the number of neurons, then  $\{i_t, f_t, c_t, o_t, h_t\} \in \mathbb{R}^{n \times 1}$ . For simplicity, in the rest of the paper to refer to the weights and biases in the LSTM model, we use column vectors  $w_{lstm}$  and  $b_{lstm}$  which include all the elements in  $\{W_{xi}, W_{ci}, W_{xf}, W_{hf}, W_{cf}, W_{xc}, W_{hc}, W_{xo}, W_{ho}, W_{co}\}$  and  $\{b_i, b_f, b_c, b_o\}$  respectively. The LSTM equations can be briefly written as follows

$$\begin{cases} c_t = f(c_{t-1}, h_{t-1}, x_t; w_{lstm}, b_{lstm}) \\ h_t = g(h_{t-1}, c_{t-1}, x_t; w_{lstm}, b_{lstm}) \end{cases} \quad (6)$$

where  $g(\cdot)$  and  $f(\cdot)$  can be derived from Eqs. (1) to (5).

## 3. Training LSTM for time-series prediction

Given a uni-variate or a multi-variate time-series, one may consider the whole time-series as a sample and train the LSTM network. For systems with complex dynamics, such as chaotic

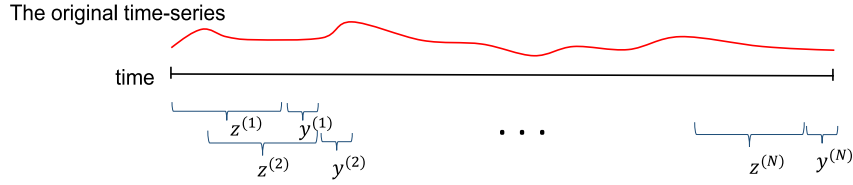


Fig. 2. Sampling from the original time-series.

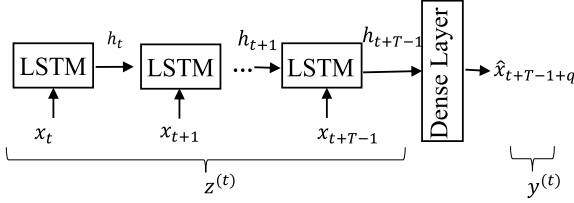


Fig. 3. Unfolded LSTM model.

systems, Suykens, Vandewalle, and De Moor (1996) showed how one can improve trajectory learning by training on packets of increasing size. In order to improve the learning process and performance, we define smaller sub-samples from the original time-series. As shown in Fig. 2, we deploy a moving window to sample from the original time-series. Note that the illustrations are done for a uni-variate time-series, and it can be extended to a multi-variate case.

Considering  $T$  as the sequence length, the sequence  $z^{(t)} = [x_t; x_{t+1}; \dots; x_{t+T-1}] \in \mathbb{R}^T$  and  $y^{(t)} = x_{t+T-1+q} \in \mathbb{R}$  are the  $t$ th input and output of the LSTM network, respectively. Note that  $q$  is a positive integer that indicates the number of steps ahead to be predicted. Also,  $N$  is the total number of subsamples and it depends on the length of the original time-series and the sequence length  $T$ .

Fig. 3 depicts how LSTM models process sequential data. The values of gates, memory cell and hidden states are updated in a recurrent way. As shown in the unfolded scheme, in each time step, the input of the LSTM cell includes the input of the corresponding time step and the previous hidden state values and the hidden values of the last LSTM cell are considered as the input of the dense layer.

To predict the future behavior of the system, one may use a dense layer as follows

$$\hat{y}^{(t)} = w_{\text{dense}}^T h_{t+T-1} + b_{\text{dense}}, \quad t = 1, \dots, N, \quad (7)$$

where  $w_{\text{dense}} \in \mathbb{R}^{n \times 1}$  and  $b_{\text{dense}} \in \mathbb{R}$  are the weights and bias term in the dense layer, and  $h_{t+T-1} \in \mathbb{R}^{n \times 1}$  represents the hidden state of the last LSTM cell. Note that the dimension of the hidden state is equal to the predefined number of neurons. Clearly, in the last step, the hidden state values include the information of all previous inputs as it is a function of the cell memory and the output gate as shown in (5).

In this paper, we use a quadratic loss function to train the network for the regression problem and deploy  $L_2$ -norm regularization to avoid overfitting. By denoting  $w$  and  $b$  all the parameters of the LSTM layer ( $w_{\text{lstm}}$  and  $b_{\text{lstm}}$ ) and the dense layer ( $w_{\text{dense}}$  and  $b_{\text{dense}}$ ), the objective function is defined as follows

$$(\hat{w}_{\text{lstm}}, \hat{w}_{\text{dense}}, \hat{b}_{\text{lstm}}, \hat{b}_{\text{dense}}) = (\hat{w}, \hat{b}) = \arg \min_{w, b} J \quad (8)$$

$$J = \frac{1}{N} \sum_{t=1}^N (\hat{y}^{(t)} - y^{(t)})^2 + \gamma w^T w$$

where  $\gamma \in \mathbb{R}^+$  is the regularization parameter that has to be tuned and  $\hat{y}^{(t)}$  is the estimated value for the output of the system

when  $z^{(t)}$  is the input sequence. Given an unseen data point  $z^{(\eta)}$  with length  $T$  and  $z_{t'}^{(\eta)}$  representing the  $t'$ th element in the sequence, the hidden state values can be evaluated as follows

$$h_{t'} = g(h_{t'-1}, c_{t'-1}, z_{t'}^{(\eta)}; \hat{w}_{\text{lstm}}, \hat{b}_{\text{lstm}}), \quad (9)$$

where  $t' = \eta, \dots, \eta+T-1$ . Then, the prediction is done as follows

$$\hat{y}^{(\eta)} = \hat{w}_{\text{dense}}^T h_{\eta+T-1} + \hat{b}_{\text{dense}}. \quad (10)$$

#### 4. Transductive LSTM for time series prediction

##### 4.1. Optimization problem

In this section, we discuss the proposed *Transductive LSTM* (T-LSTM) as a localized version of LSTM models. As discussed in previous section, LSTM models can be used for regression problems and in the corresponding model, all subsamples have the same impact on the model parameters ( $w, b$ ). However, in the proposed T-LSTM, the influence of each training data point on the parameters depends on its similarity to the test data point, which is a new sequence with length  $T$ . Therefore, the emphasis of the training phase is to obtain a better performance in the vicinity of the test point, and higher importance is given to the performance of the model for estimating the training samples that are in the neighborhood of the test point. Hence, it can be concluded that as a result all model parameters depend on the test point.

Assuming  $z^{(\eta)}$  being an unseen sequence, the state space model of the T-LSTM can be written as follows

$$\begin{cases} c_{t,\eta} = f(c_{t-1,\eta}, h_{t-1,\eta}, x_t; w_{\text{lstm},\eta}, b_{\text{lstm},\eta}) \\ h_{t,\eta} = g(h_{t-1,\eta}, c_{t-1,\eta}, x_t; w_{\text{lstm},\eta}, b_{\text{lstm},\eta}). \end{cases} \quad (11)$$

Note that the model described in (11) is different from the one in (6): in (6) the model parameters are independent from the test point while in (11) the model parameters depend on the feature vector of the test point. We use  $\eta$  as a subscript to indicate that the model parameters depend on the new given data point  $z^{(\eta)}$ . It is important to mention that the test label is considered to be unseen and the only use of the test point in the training phase is to affect the importance of the training data points based on the similarity between their feature vector and the test point feature vector. Afterwards, the prediction using the dense layer can be written as follows

$$\hat{y}_\eta^{(t)} = w_{\text{dense},\eta}^T h_{t+T-1,\eta} + b_{\text{dense},\eta}, \quad t = 1, \dots, N, \quad (12)$$

where  $w_{\text{dense},\eta} \in \mathbb{R}^{n \times 1}$  and  $b_{\text{dense},\eta} \in \mathbb{R}$  are the weights and bias term in the dense layer.

To specify the dependencies on the new unseen data point, assume  $s_{t,\eta}$  being the similarity between length  $T$  sequences  $z^{(t)}$  and  $z^{(\eta)}$ , and  $w_\eta$  and  $b_\eta$  denoting all the parameters in ( $w_{\text{lstm},\eta}, w_{\text{dense},\eta}$ ) and ( $b_{\text{lstm},\eta}, b_{\text{dense},\eta}$ ) respectively. The objective function can be written as follows

$$(\hat{w}_{\text{lstm},\eta}, \hat{w}_{\text{dense},\eta}, \hat{b}_{\text{lstm},\eta}, \hat{b}_{\text{dense},\eta}) = (\hat{w}_\eta, \hat{b}_\eta) = \arg \min_{w_\eta, b_\eta} J_\eta \quad (13)$$

$$J_\eta = \frac{1}{N} \sum_{t=1}^N s_{t,\eta} (\hat{y}_\eta^{(t)} - y^{(t)})^2 + \gamma w_\eta^T w_\eta$$



where  $s_{t,\eta} \in \mathbb{R}^+$ . Note that the main difference between (8) and (13) is the fact that in the latter, the model parameters depend on the new point  $z^{(\eta)}$ . Thus, for each unseen subsample the model has to be retrained. This means that all the parameters  $\hat{w}_\eta$  and  $\hat{b}_\eta$  can be different for each test point. Therefore, although the time complexity of the proposed method is similar to the one of inductive LSTM, the bottleneck of the transductive one is the number of models that have to be trained. This number depends on the number of data points in the test set.

The similarity function  $s_{t,\eta}$  can take binary or positive real values. In case of a binary value,  $k$ -Nearest Neighbors of  $z^{(\eta)}$  are selected as training data and the rest of the samples are discarded. In this case,  $s_{t,\eta}$  for all selected samples are equal to one. Hence, their impact on the training phase is the same. In addition, one may deploy a clustering approach to select the samples similar to the test point. Bandara et al. (2017) utilized clustering to group the data and then trained a separate LSTM model per cluster. Note that in this case the number of data points to train the model is reduced.

On the other hand, in case  $s_{t,\eta}$  are real-valued, all subsamples are taken into account in the training phase, while their impact on the training phase is different. In this case  $s_{t,\eta}$  can be calculated using a similarity function, e.g. Gaussian or cosine similarity. In this study we deploy the real-valued  $s_{t,\eta}$  and use cosine similarity to weight the training samples. Cosine similarity is defined as  $\kappa_{t,\eta} = \frac{z^{(t)T} z^{(\eta)}}{\|z^{(t)}\| \|z^{(\eta)}\|}$  and it takes a value between  $-1$  and  $1$ . To achieve a non-negative similarity value, we perform two transformations on the cosine similarity: (1) linear transformation at which point  $s_{t,\eta} = \kappa_{t,\eta} + 1$  such that the similarity is between  $0$  and  $2$ , and (2) non-linear sigmoid transformation at which point  $s_{t,\eta} = \sigma(\kappa_{t,\eta}) = \frac{1}{1 + \exp(-\kappa_{t,\eta})}$  such that the similarity is between  $0$  and  $1$ . Note that, other similarity functions such as the Gaussian similarity can be used to specify the weights; however, in this study we deploy cosine similarity to avoid additional tuning parameters.

Note that using T-LSTM is computationally expensive; so, one may consider if this method is suitable for modeling based on the characteristics of the problem. In inductive approach, the parameters are calculated regardless of the test point, while in the transductive one separate model is trained for each test point. Thus, T-LSTM is not suitable for problem with large test sets as the model has to be trained for each test point separately. It can be considered suitable method for time-series prediction where there is enough time for retraining the model. In addition, T-LSTM is suitable for cases where the relation between variables in different parts of the input space is different. For example, possible differences in patterns in summer and winter is our intuition to use a transductive approach in weather forecasting.

For a given point  $z^{(\eta)}$ , the updates of the hidden state are done as follows

$$h_{t',\eta} = g(h_{t'-1,\eta}, c_{t'-1,\eta}, z_{t'}^{(\eta)}; \hat{w}_{\text{lstm},\eta}, \hat{b}_{\text{lstm},\eta}) \quad (14)$$

where  $t' = \eta, \dots, \eta + T - 1$ . Afterwards, the final prediction can be obtained as follows

$$\hat{y}_\eta^{(\eta)} = \hat{w}_{\text{dense},\eta}^T h_{\eta+T-1,\eta} + \hat{b}_{\text{dense},\eta}. \quad (15)$$

Note that (9), (10), (14) and (15) refer to the prediction of the system when  $z^{(\eta)}$  is the input. However, the hidden state updates and the prediction model in the T-LSTM depends also implicitly on the new point  $z^{(\eta)}$  feature vector as the model parameters are optimized based on the similarity between the training points and  $z^{(\eta)}$ . This is indicated by  $\eta$  subscript in  $h_{t',\eta}$  and  $\hat{y}_\eta^{(\eta)}$  in (14) and (15).



Fig. 4. Weather stations (picture from Google maps).

## 4.2. Tuning parameters

Similar to (8),  $\gamma_\eta$  is a tuning parameter in (13). In addition, the number of neurons in LSTM gates is a tuning parameter both in inductive and transductive approaches.

Since the emphasis of the T-LSTM is to achieve a good performance in the vicinity of the test point, one may select the most similar samples to the test point as the validation set. However, in this scenario the model cannot exploit the information of most similar data points in the training phase. In this study, we assume that there is no sudden change in the pattern of the original series; thus, the distance between two consecutive samples in the feature space is not large. Therefore, the last samples of the training period, which include the samples prior to the test point, should be in the neighborhood of the test point in the feature space. Hence, these samples are selected as the validation set.

## 5. Experiments

### 5.1. Dataset

In this study, data have been collected from the Weather Underground website and include real measurements of weather variables such as minimum and maximum temperature, dew point, and wind speed for 5 cities including Brussels, Antwerp, Liege, Amsterdam and Eindhoven, as shown in Fig. 4.

The data cover a time period from the beginning of 2007 to mid-2014 and contain 18 measured weather variables for each day per city. To exploit all available information, the number of training data points is different for each test set, and it is equal to the number of days from the beginning of 2007 until the day before the test set.

To evaluate the performance of the proposed methods in various weather conditions, two test sets are defined: (i) from mid-November 2013 to mid-December 2013 (Nov/Dec) and (ii) from mid-April 2014 to mid-May 2014 (Apr/May).

### 5.2. Model selection

As the difference between weather variables in two days in a row is not significant, in this study we use the last samples of the training period as validation set. For each test set, two months prior to the test set is considered as the validation set. Hence, for the Nov/Dec test set, the validation set includes the period from mid-September 2013 to mid-November 2013 and for the Apr/May

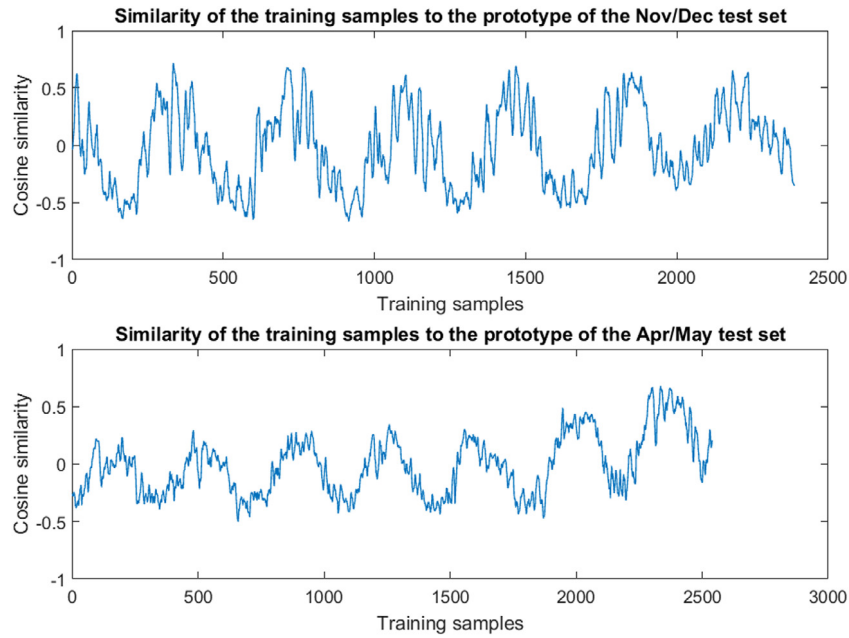


Fig. 5. Similarity of the training samples to both test sets.

test set, it includes the period from mid-February 2014 to mid-April 2014. The samples that are left for training the models are shuffled.

In this study the implementation of the LSTM models are in TensorFlow, which is based on the model described by Zaremba et al. (2014). Note that this model is slightly different from the one discussed by Graves (2013). The tuning parameters are the number of neurons and the regularization parameter. For the number of neurons in the LSTM, we try the values: {16, 32, 64, 128, 256}, and for the regularization parameter  $\gamma$  the values {0.1, 0.2, 1, 2, 10} are considered. For the inner state, mostly tanh is used as the activation function. In addition to that, in Section 5.3 we discuss the performance while sigmoid is chosen as the activation function. The ADAM optimizer is used in this paper to optimize the objective function. Note that the total number of variables in each time step is equal to  $18 \times 5$  and the number of samples depends on the sequence length  $T$ .

### 5.3. Results

As mentioned before, we use cosine-based similarity to specify the weights of the samples in the objective function. In Fig. 5 the cosine similarity of the training data points to the prototypes of both tests sets is depicted. As can be seen the similarity values have a temporal meaning; i.e. the samples in a specific period of the year have high similarity value with each other. In addition, it can be concluded that there are no significant changes in the weather variables for two days in a row (i.e. temporally close samples).

The performance of the proposed method is evaluated on an application of maximum and minimum temperature prediction in Brussels, Belgium. The forecasting is done for one to six days ahead. To avoid local minima problems in LSTMs, we did the experiments five times. Tables 1 and 2 present the median of Mean Absolute Error (MAE) and Mean Squared Error (MSE) of the inductive and transductive approaches on the test sets for different internal activation functions and sequence lengths. The bold one shows the best performance for the corresponding setting. The results for the transductive approaches are presented for both linear ( $\kappa_{t,\eta} + 1$ ) and non-linear ( $\sigma(\kappa_{t,\eta})$ ) transformation of the cosine similarity.

Note that for  $i$ th day in the test periods, the model is using the actual value (and not the predicted value of previous days). For example, assuming the data is available until 14th of April, then model for one day ahead predicts 15th April, for second days ahead predicts 16th April and so on. As one day is passed, the actual value is used. So, the actual value of the 15th April is used as input feature to predict 16th April (one day ahead).

As is shown, the transductive approaches outperform the inductive approach in the case that the transfer function is a sigmoid. In case of a hyperbolic tangent transfer function, there are a few cases that the inductive approach has better performance than the transductive approach. In addition, there is no significant difference between the linear and non-linear transformation of the cosine distance.

Fig. 6 presents the statistical distribution of MAE for both test sets together for inductive LSTM and T-LSTM. It is clear that T-LSTM shows better performance than LSTM. This is due to the fact that the relation between weather variables can be different in different time of a year.

In order to be more specific, Figs. 7 and 8 present the statistical distribution of MAE for Nov/Dec test set and Apr/May test set for inductive LSTM and T-LSTM. It can be seen that the performances of T-LSTM and LSTM are relatively competitive for Nov/Dec while in Apr/May test set, T-LSTM outperforms LSTM. Hence, it can be concluded that in Nov/Dec the inductive approach fits a model that works well in Nov/Dec period; i.e. for this period, the transductive approach cannot improve the performance. On the contrary, for Apr/May test set, the model fitted by inductive approach can be improved by considering different impact for training data points.

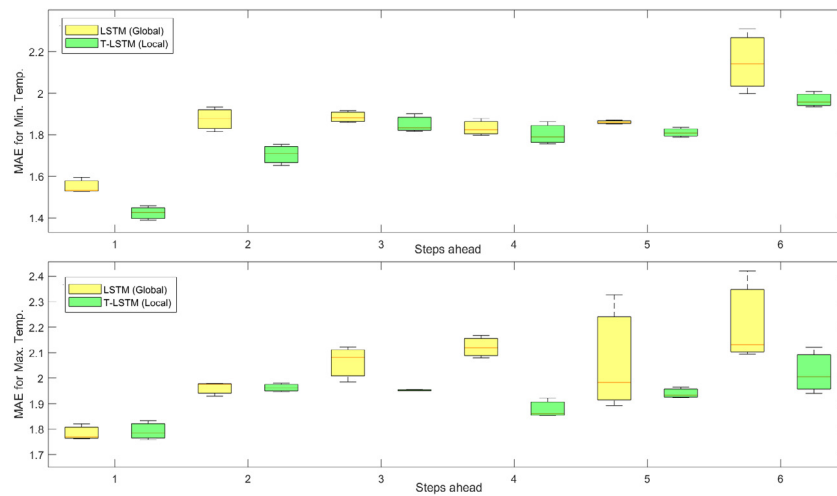
In Figs. 9 and 10, the performance of the LSTM and T-LSTM models, as data-driven approaches, are compared with the performance of ridge regression, as a baseline regularization method, and Weather Underground company ([www.wunderground.com](http://www.wunderground.com)), as a state-of-the-art method for weather forecasting. Similar to LSTM and T-LSTM, the number of delays and regularization parameter are tuned in ridge regression method. The results reveal that while the information that is used in the data-driven approaches is very limited (only few cities and few years of data

**Table 1**  
MAE and MSE of minimum and maximum temperature prediction for LSTM, T-LSTM with both linear and non-linear transformation of cosine similarity, for sequence length 10 and 20, and for 1 to 6 days ahead for Nov/Dec test set.

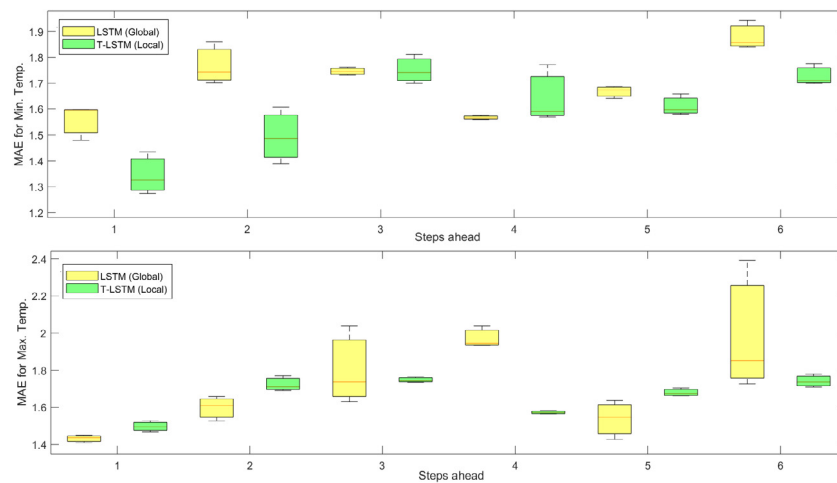
Seq. Len.	Steps ahead	Temp.	Activation function: <i>Tanh</i>						Activation function: <i>Sigmoid</i>					
			Mean absolute error			Mean squared error			Mean absolute error			Mean squared error		
			LSTM	T-LSTM ( $\kappa + 1$ )	T-LSTM ( $\sigma(\kappa)$ )	LSTM	T-LSTM ( $\kappa + 1$ )	T-LSTM ( $\sigma(\kappa)$ )	LSTM	T-LSTM ( $\kappa + 1$ )	T-LSTM ( $\sigma(\kappa)$ )	LSTM	T-LSTM ( $\kappa + 1$ )	T-LSTM ( $\sigma(\kappa)$ )
10	1	Min	1.54	1.54	<b>1.50</b>	3.84	<b>3.74</b>	<b>3.74</b>	1.72	<b>1.61</b>	1.66	4.61	<b>4.10</b>	4.20
		Max	1.88	1.32	<b>1.30</b>	4.58	2.92	<b>2.82</b>	1.52	1.41	<b>1.37</b>	3.80	<b>3.49</b>	3.52
	2	Min	<b>1.74</b>	1.83	<b>1.74</b>	<b>4.32</b>	4.70	4.42	1.89	<b>1.75</b>	1.76	5.48	<b>4.37</b>	5.35
		Max	2.16	<b>1.88</b>	1.94	6.47	5.13	<b>5.10</b>	1.61	<b>1.55</b>	1.59	4.41	<b>3.87</b>	3.99
	3	Min	<b>1.74</b>	1.77	<b>1.74</b>	4.40	4.46	<b>4.35</b>	2.08	1.95	<b>1.91</b>	6.97	5.35	<b>5.29</b>
		Max	<b>1.95</b>	1.99	<b>1.95</b>	<b>4.73</b>	4.95	5.19	2.10	<b>1.79</b>	1.81	7.51	6.26	<b>5.90</b>
	4	Min	<b>1.56</b>	1.67	1.57	<b>3.73</b>	4.31	3.80	2.02	1.85	<b>1.70</b>	6.36	5.50	<b>4.60</b>
		Max	1.94	<b>1.59</b>	1.86	4.45	<b>3.96</b>	4.34	1.97	<b>1.90</b>	1.92	6.46	<b>5.79</b>	5.89
	5	Min	<b>1.67</b>	1.81	<b>1.67</b>	4.47	4.90	<b>4.30</b>	2.01	1.92	<b>1.76</b>	5.84	5.40	<b>4.47</b>
		Max	1.95	1.86	<b>1.59</b>	5.44	4.42	<b>3.93</b>	1.89	<b>1.81</b>	1.91	5.79	<b>5.40</b>	5.99
	6	Min	1.96	2.22	<b>1.87</b>	5.81	<b>5.18</b>	5.25	2.02	1.90	<b>1.73</b>	5.50	4.80	<b>4.28</b>
		Max	1.90	1.83	<b>1.76</b>	5.98	5.90	<b>5.13</b>	2.09	<b>1.73</b>	1.83	6.96	<b>5.24</b>	5.49
20	1	Min	1.76	1.58	<b>1.49</b>	5.40	4.07	<b>3.67</b>	1.74	<b>1.65</b>	1.69	4.53	<b>4.19</b>	4.40
		Max	1.48	1.39	<b>1.34</b>	3.32	3.10	<b>2.85</b>	1.43	1.43	<b>1.35</b>	3.49	3.59	<b>3.38</b>
	2	Min	1.77	1.80	<b>1.71</b>	4.40	4.51	<b>4.28</b>	1.85	<b>1.74</b>	<b>1.74</b>	5.63	5.06	<b>4.92</b>
		Max	1.83	1.74	<b>1.51</b>	4.55	<b>4.09</b>	4.26	1.60	<b>1.48</b>	1.49	4.15	<b>3.53</b>	3.62
	3	Min	1.77	1.80	<b>1.75</b>	<b>4.46</b>	4.61	4.60	2.03	1.96	<b>1.94</b>	6.26	5.55	<b>5.43</b>
		Max	<b>1.73</b>	2.09	1.84	<b>3.90</b>	5.60	4.51	1.92	<b>1.73</b>	1.76	5.67	5.76	<b>5.57</b>
	4	Min	1.62	1.62	<b>1.60</b>	3.92	<b>3.68</b>	3.86	1.97	<b>1.69</b>	1.72	6.04	<b>4.50</b>	4.74
		Max	1.99	1.99	<b>1.91</b>	4.70	5.04	<b>4.43</b>	1.94	1.86	<b>1.80</b>	4.98	5.10	<b>4.65</b>
	5	Min	<b>1.68</b>	<b>1.68</b>	1.69	4.48	4.38	<b>4.31</b>	2.00	<b>1.68</b>	1.76	5.75	<b>4.81</b>	4.78
		Max	<b>1.54</b>	2.09	1.57	<b>3.37</b>	5.65	3.68	2.01	<b>1.75</b>	1.76	4.99	<b>4.25</b>	4.30
	6	Min	<b>1.85</b>	1.87	1.90	<b>5.28</b>	5.31	5.42	1.98	<b>1.76</b>	1.77	5.30	4.46	<b>4.28</b>
		Max	1.85	1.79	<b>1.76</b>	6.94	6.26	<b>5.46</b>	1.94	1.77	<b>1.71</b>	6.19	5.40	<b>5.22</b>

**Table 2**  
MAE and MSE of minimum and maximum temperature prediction for LSTM, T-LSTM with both linear and non-linear transformation of cosine similarity, for sequence length 10 and 20, and for 1 to 6 days ahead for Apr/May test set.

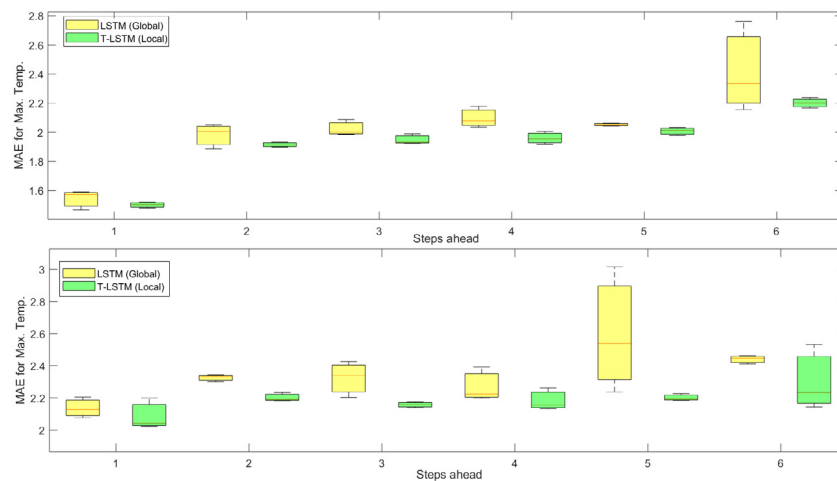
Seq. Len.	Steps ahead	Temp.	Activation function: <i>Tanh</i>						Activation function: <i>Sigmoid</i>					
			Mean absolute error			Mean squared error			Mean absolute error			Mean squared error		
			LSTM	T-LSTM ( $\kappa + 1$ )	T-LSTM ( $\sigma(\kappa)$ )	LSTM	T-LSTM ( $\kappa + 1$ )	T-LSTM ( $\sigma(\kappa)$ )	LSTM	T-LSTM ( $\kappa + 1$ )	T-LSTM ( $\sigma(\kappa)$ )	LSTM	T-LSTM ( $\kappa + 1$ )	T-LSTM ( $\sigma(\kappa)$ )
10	1	Min	1.74	<b>1.52</b>	1.68	4.08	<b>3.94</b>	4.41	1.63	1.59	<b>1.56</b>	4.04	<b>3.73</b>	3.88
		Max	2.15	2.10	<b>2.04</b>	<b>7.44</b>	6.78	<b>7.44</b>	2.21	<b>2.12</b>	2.16	7.14	<b>7.07</b>	7.36
	2	Min	2.05	1.99	<b>1.95</b>	<b>7.42</b>	7.69	8.86	2.33	1.96	<b>1.91</b>	7.95	7.30	<b>7.02</b>
		Max	2.33	2.23	<b>2.19</b>	8.24	8.16	<b>8.09</b>	2.68	<b>2.33</b>	2.47	9.91	<b>8.23</b>	8.74
	3	Min	2.07	<b>1.98</b>	2.02	8.54	7.59	<b>7.54</b>	2.08	1.96	<b>1.95</b>	7.22	<b>7.07</b>	7.22
		Max	2.71	<b>2.15</b>	2.17	12.77	<b>7.30</b>	7.50	2.55	2.37	<b>2.36</b>	9.10	<b>7.89</b>	8.20
	4	Min	2.16	<b>1.93</b>	2.14	8.47	<b>6.42</b>	6.60	2.07	<b>1.92</b>	2.19	6.63	<b>6.29</b>	7.04
		Max	2.87	<b>2.22</b>	<b>2.22</b>	13.41	7.97	<b>7.86</b>	2.36	2.30	<b>2.25</b>	9.04	8.65	<b>7.96</b>
	5	Min	2.05	2.02	<b>2.01</b>	<b>6.91</b>	7.17	7.30	2.35	<b>2.11</b>	2.22	7.84	<b>6.94</b>	7.18
		Max	2.85	2.52	<b>2.26</b>	13.31	10.59	<b>8.32</b>	2.53	2.28	<b>2.19</b>	10.59	8.86	<b>8.33</b>
	6	Min	2.66	<b>2.20</b>	2.28	13.80	8.35	<b>7.48</b>	2.52	<b>2.35</b>	2.54	8.58	<b>7.95</b>	8.61
		Max	2.95	<b>2.58</b>	2.72	14.40	<b>11.38</b>	12.70	2.44	<b>2.36</b>	2.43	10.38	<b>8.54</b>	10.38
20	1	Min	1.62	<b>1.52</b>	1.58	4.25	<b>4.08</b>	4.42	1.57	<b>1.50</b>	1.67	4.03	<b>3.80</b>	4.03
		Max	2.20	<b>2.10</b>	2.13	7.60	<b>7.03</b>	7.33	2.16	<b>2.11</b>	2.13	7.56	<b>7.43</b>	7.89
	2	Min	2.00	1.92	<b>1.91</b>	7.54	7.90	<b>7.17</b>	2.03	1.94	<b>1.90</b>	7.95	<b>7.38</b>	7.56
		Max	2.38	2.31	<b>2.25</b>	8.33	<b>7.81</b>	8.00	2.61	2.57	<b>2.50</b>	9.66	9.36	<b>8.99</b>
	3	Min	2.17	<b>1.93</b>	2.05	8.32	<b>6.21</b>	7.54	2.00	<b>1.97</b>	1.98	8.39	7.77	<b>7.39</b>
		Max	2.66	2.39	<b>2.20</b>	9.50	8.85	<b>7.37</b>	2.33	<b>2.23</b>	2.33	7.72	<b>7.30</b>	7.73
	4	Min	2.21	2.14	<b>2.06</b>	7.58	7.05	<b>6.37</b>	2.18	2.02	<b>1.95</b>	7.10	<b>6.41</b>	6.92
		Max	2.88	<b>2.39</b>	2.49	13.44	<b>8.85</b>	8.88	2.22	<b>2.15</b>	2.20	7.48	<b>7.35</b>	7.56
	5	Min	2.37	2.23	<b>2.18</b>	9.70	9.40	<b>6.94</b>	2.28	<b>2.05</b>	2.32	7.53	<b>7.32</b>	7.65
		Max	2.94	<b>2.52</b>	2.78	12.26	10.57	<b>10.40</b>	2.73	2.54	<b>2.30</b>	12.49	9.01	<b>8.64</b>
	6	Min	2.37	<b>2.28</b>	<b>2.28</b>	9.65	<b>8.10</b>	9.70	2.45	<b>2.31</b>	2.43	8.28	<b>7.84</b>	8.15
		Max	2.76	2.67	<b>2.62</b>	10.61	<b>9.79</b>	9.82	2.79	2.51	<b>2.23</b>	11.12	10.73	<b>8.30</b>



**Fig. 6.** Comparing statistical distribution MAE of minimum and maximum temperature prediction for LSTM and T-LSTM for 1 to 6 days ahead for both test sets together.



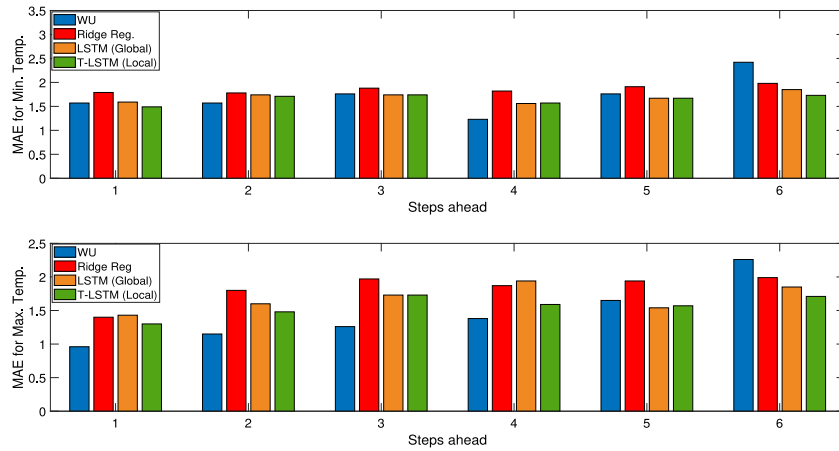
**Fig. 7.** Comparing statistical distribution MAE of minimum and maximum temperature prediction for LSTM and T-LSTM for 1 to 6 days ahead for Nov/Dec test set.



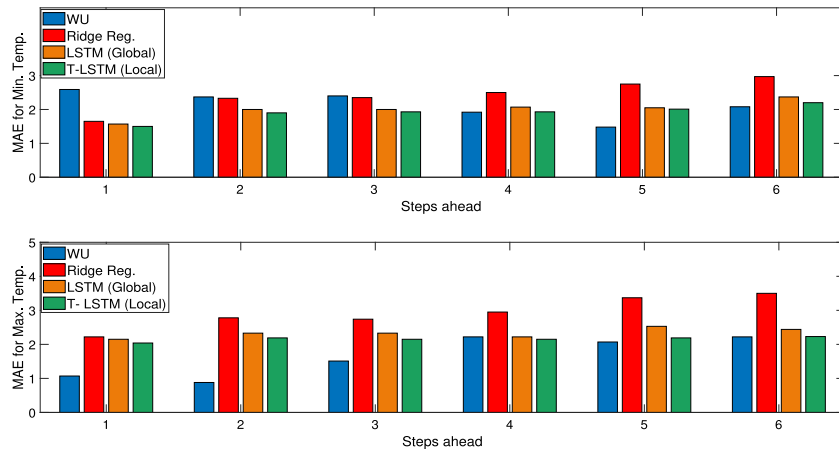
**Fig. 8.** Comparing statistical distribution MAE of minimum and maximum temperature prediction for LSTM and T-LSTM for 1 to 6 days ahead for Apr/May test set.

are taken into account), the performance is competitive with the state-of-the-art models used by Weather Underground. Moreover, it can be seen that the LSTM and T-LSTM models outperform

the baseline data-driven method. This difference between the performance of LSTM and T-LSTM methods, and the baseline method is larger for long-term prediction.



**Fig. 9.** Comparing MAE of minimum and maximum temperature prediction for Weather Underground, ridge regression, LSTM, and T-LSTM for 1 to 6 days ahead for Nov/Dec test set.



**Fig. 10.** Comparing MAE of minimum and maximum temperature prediction for Weather Underground, ridge regression, LSTM, and T-LSTM for 1 to 6 days ahead for Apr/May test set.

In addition, the comparison between the MAE of the models of inductive and transductive approaches is made in Figs. 9 and 10. In many cases the transductive approach outperforms the inductive approach or provides equally good performance. In the Apr/May test set in all cases the transductive approach shows better performance than the global one, whereas in the Nov/Dec test set the inductive and the transductive methods are more competitive. Note that the performance of inductive LSTM model in Nov/Dec is better than its performance in Apr/May.

## 6. Conclusion

In this study, we deployed an LSTM model for an application of temperature prediction. We proposed the T-LSTM model to exploit transductive learning in time-series prediction. In this case, the model can be achieved by altering the cost function of the LSTM model such that the samples in the neighborhood of the test point have a higher impact in the objective function. We tested two weighting schemes which are based on the cosine similarity between the training data points and the test point. The experiments are done for two different periods in a year to evaluate the performance of the proposed method in different weather conditions. We examined different sequence length and transfer functions and the results suggest that transductive LSTM outperforms LSTM in many cases. Moreover, despite of using

limited data, it is shown that the performance of the inductive and transductive LSTM models as data-driven approaches are competitive with the state-of-the-art methods in weather forecasting.

## Acknowledgments

Research supported by Research Council KUL, Belgium: CoE PFV/10/002 (OPTEC), PhD/Postdoc grants Flemish Government; FWO, Belgium: projects: G0A4917N (Deep restricted kernel machines), G.088114N (Tensor based data similarity), ERC Advanced Grant E-DUALITY (787960).

## References

- Bandara, K., Bergmeir, C., & Smyl, S. (2017). Forecasting across time series databases using long short-term memory networks on groups of similar series. *arXiv preprint arXiv:1710.03222*.
- Bauer, P., Thorpe, A., & Brunet, G. (2015). The quiet revolution of numerical weather prediction. *Nature*, 525(7567), 47–55.
- Bengio, Y., Simard, P., & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2), 157–166.
- Bottou, L., & Vapnik, V. (1992). Local learning algorithms. *Neural Computation*, 4(6), 888–900.



- Chen, L., Zhou, M., Wu, M., She, J., Liu, Z., Dong, F., et al. (2018). Three-layer weighted fuzzy support vector regression for emotional intention understanding in human-robot interaction. *IEEE Transactions on Fuzzy Systems*.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., et al. (2014). Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *EMNLP* (pp. 1724–1734).
- Do, T.-N., & Poulet, F. (2015). Random local SVMs for classifying large datasets. In *International conference on future data and security engineering* (pp. 3–15). Springer.
- Elattar, E. E., Goulermas, J., & Wu, Q. H. (2010). Electric load forecasting based on locally weighted support vector regression. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 40(4), 438–447.
- Elman, J. L. (1990). Finding structure in time. *Cognitive Science*, 14(2), 179–211.
- Freeman, B. S., Taylor, G., Gharabaghi, B., & Thé, J. (2018). Forecasting air quality time series using deep learning. *Journal of the Air & Waste Management Association*.
- Gers, F. A., & Schmidhuber, J. (2000). Recurrent nets that time and count. In *IJCNN, Vol. 3* (pp. 189–194). IEEE.
- Gers, F. A., Schraudolph, N. N., & Schmidhuber, J. (2002). Learning precise timing with LSTM recurrent networks. *Journal of Machine Learning Research (JMLR)*, 3(Aug), 115–143.
- Gilardi, N., & Bengio, S. (2000). Local machine learning models for spatial data analysis. *Journal of Geographic Information and Decision Analysis*, 4, 11–28.
- Graves, A. (2013). Generating sequences with recurrent neural networks. arXiv preprint arXiv:1308.0850.
- Graves, A., Mohamed, A.-R., & Hinton, G. (2013). Speech recognition with deep recurrent neural networks. In *ICASSP* (pp. 6645–6649). IEEE.
- Graves, A., & Schmidhuber, J. (2005). Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks*, 18(5–6), 602–610.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.
- Karevan, Z., Feng, Y., & Suykens, J. A. K. (2016). Moving least squares support vector machines for weather temperature prediction. In *ESANN* (pp. 611–616).
- Karevan, Z., & Suykens, J. A. K. (2018). Transductive feature selection using clustering-based sample entropy for temperature prediction in weather forecasting. *Entropy*, 20(4), 264.
- Levin, D. (1998). The approximation power of moving least-squares. *Mathematics of Computation of the American Mathematical Society*, 67(224), 1517–1531.
- Lipton, Z. C., Kale, D. C., Elkan, C., & Wetzell, R. (2016). Learning to diagnose with LSTM recurrent neural networks. *ICLR*.
- Ng, J. Y.-H., Hausknecht, M., Vijayanarasimhan, S., Vinyals, O., Monga, R., & Toderici, G. (2015). Beyond short snippets: Deep networks for video classification. In *CVPR, 2015* (pp. 4694–4702). IEEE.
- Pang, S., & Kasabov, N. (2004). Inductive vs transductive inference, global vs local models: SVM, TSVM, and SVMT for gene expression classification problems. In *Neural networks, 2004. proceedings. 2004 IEEE international joint conference on, Vol. 2* (pp. 1197–1202). IEEE.
- Schmidhuber, J., Wierstra, D., Gagliolo, M., & Gomez, F. (2007). Training recurrent networks by evoluno. *Neural Computation*, 19(3), 757–779.
- Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Advances in neural information processing systems* (pp. 3104–3112).
- Suykens, J. A. K., Vandewalle, J. P. L., & De Moor, B. L. (1996). *Artificial neural networks for modelling and control of non-linear systems*. Springer Science & Business Media.
- Tian, Y., & Pan, L. (2015). Predicting short-term traffic flow by long short-term memory recurrent neural network. In *Smart city/socialcom/sustaincom (smartcity), 2015 IEEE international conference on* (pp. 153–158). IEEE.
- Vapnik, V. (1992). Principles of risk minimization for learning theory. In *Advances in neural information processing systems* (pp. 831–838).
- Wu, Q., Li, H., Meng, F., & Ngan, K. N. (2018). Generic proposal evaluator: A lazy learning strategy toward blind proposal quality assessment. *IEEE Transactions on Intelligent Transportation Systems*, 19(1), 306–319.
- Zaremba, W., Sutskever, I., & Vinyals, O. (2014). Recurrent neural network regularization. arXiv preprint arXiv:1409.2329.
- Zhai, D., Liu, X., Chang, H., Zhen, Y., Chen, X., Guo, M., et al. (2018). Parametric local multiview Hamming distance metric learning. *Pattern Recognition*, 75, 250–262.