# Temporal Logistic Neural Bag-of-Features for Financial Time series Forecasting leveraging Limit Order Book Data

Nikolaos Passalis, Anastasios Tefas, Juho Kanniainen, Moncef Gabbouj
and Alexandros Iosifidis

## Abstract

Time series forecasting is a crucial component of many important applications, ranging from forecasting the stock markets to energy load prediction. The high-dimensionality, velocity and variety of the data collected in these applications pose significant and unique challenges that must be carefully addressed for each of them. In this work, a novel Temporal Logistic Neural Bag-of-Features approach, that can be used to tackle these challenges, is proposed. The proposed method can be effectively combined with deep neural networks, leading to powerful deep learning models for time series analysis. However, combining existing BoF formulations with deep feature extractors pose significant challenges: the distribution of the input features is not stationary, tuning the hyper-parameters of the model can be especially difficult and the normalizations involved in the BoF model can cause significant instabilities during the training process. The proposed method is capable of overcoming these limitations by a employing a novel adaptive scaling mechanism and replacing the classical Gaussian-based density estimation involved in the regular BoF model with a logistic kernel. The effectiveness of the proposed approach is demonstrated using extensive experiments on a large-scale financial time series dataset that consists of more than 4 million limit orders.

*Nikolaos Passalis, Juho Kanniainen and Moncef Gabbouj are with the Faculty of Information Technology and Communication, Tampere University, Finland. Anastasios Tefas is with the School of Informatics, Aristotle University of Thessaloniki, Greece. Alexandros Iosifidis is with the Department of Engineering, Electrical and Computer Engineering, Aarhus University, Denmark. E-mail: nikolaos.passalis@tuni.fi, tefas@csd.auth.gr, juho.kanniainen@tuni.fi, moncef.gabbouj@tuni.fi, alexandros.iosifidis@eng.au.dk

# 1 Introduction

Time series forecasting is a crucial component of many important applications, ranging from predicting the behavior of financial markets [5], to accurate energy load prediction [13]. Even though the large amount of data that can be nowadays collected from these domains provide an unprecedented opportunity for applying powerful deep learning (DL) methods [23, 41, 24], the high-dimensionality, velocity and variety of such data also pose significant and unique challenges that must be carefully addressed for each application. To this end, many methods have been proposed to analyze and forecast time series data. For example, traditional approaches employ adaptive distance metrics, such as Dynamic Time Wrapping [4], to tackle these kind of tasks. However, with the advent of DL the interest is gradually shifting toward using neural network-based methods, including recurrent and convolutional architectures [25, 7], that seem to be more effective for handling such kind of data. It is worth noting that other approaches for time series analysis also exist, such as using the Bag-of-Features model (BoF) [35]. The BoF model was recently adapted toward efficiently processing large amounts of complex and high-dimensional time series [2, 1, 32], due its ability to analyze objects that consist of a varying number of features, as well as withstanding distribution shifts better than competitive methods [29].

The Bag-of-Features model (BoF) involves the following pipeline [35]: a) Several feature vectors are extracted from each input object, e.g., an image or time series. This step is called *feature extraction* and allows for forming the *feature space*, where each object is represented as a set of feature vectors. b) A set of representative feature vectors (also called *codewords*) are learned and used to quantize the extracted feature vectors. This step is called *dictionary learning*, while the learned codewords form the *dictionary* (which also called *codebook*). c) The quantized feature vectors are aggregated in order to extract a constant length representation that describes the semantic content/temporal behavior/etc. of each input object.

The BoF model is able to successfully handle objects of various lengths, providing an important advantage over other methods, since it allows for efficiently extracting a constant length representation of time series regardless their actual length. Indeed, the ability of the BoF-based model to tackle several time series analysis tasks have been demonstrated in the literature [2, 1, 32]. However, these approaches mainly employ shallow models that use simple hand-crafted features, instead of using more powerful deep feature extraction layers for extracting higher level features

2

that can better model the dynamics of a time series [7, 39]. In this work we argue that combining the BoF model with such architectures can significantly improve the performance of time series forecasting algorithms, since the BoF model allows for dealing with time series of arbitrary lengths and withstanding mild distribution shifts, while using deep feature extractors, such as recurrent and convolutional layers, allows for taking into account the more detailed temporal dynamics.

The main contribution of this work is the proposal of a novel logistic formulation of the Bag-of-Features model that is adapted toward the needs of time series forecasting and can be effectively combined with deep neural networks. The proposed method is indeed capable of combining the advantages of the BoF model with the enormous learning capacity of deep learning models, allowing for developing powerful forecasting models. However, combining existing BoF formulations with deep feature extractors pose significant challenges: the distribution of the input features is not stationary, tuning the hyper-parameters of the model can be especially difficult and the normalizations involved in the BoF formulation can cause significant instabilities during the training process. The later was found to be the main cause for the difficulties in training deep models that employ BoF layers and it is addressed by proposing an appropriate adaptive scaling approach. Furthermore, the classical Gaussian-based density estimation involved in the regular BoF model is replaced by a logistic kernel, following a probabilistic formulation of the BoF model [3], allowing for further improving the performance of the model and simplifying the implementation without requiring any sophisticated initialization scheme or careful finetuning of any hyper-parameter. Furthermore, the proposed method is able to perform fine-grained temporal modeling, as shown in Fig. 1, where the short-term, mid-term, and long-term behavior of time series are modeled. The proposed method is extensively evaluated using a large-scale financial time series dataset that consists of more than 4 million limit orders.

The rest of the paper is structured as follows. First, the related work is briefly introduced and compared to the proposed approach in Section 2. Then, the proposed method is introduced in Section 3, while the experimental evaluation is provided in Section 4. Finally, conclusions are drawn in Section 5.

# 2 Related Work

This work is mainly related to time series analysis using the BoF model. An increasing number of recent works employ variants of the Bag-of-Features model to perform time series analysis, e.g., forecasting, retrieval, etc. In [16], a BoF-based method was proposed for extracting discriminative representations by employing a discriminative objective for optimizing the codebook. A dictionary learning methods for the BoF model was also utilized in [28], in order to learn retrieval-oriented representations. A discriminant BoF approach for learning representations for action recognition was proposed in [17], while a dynemes-based one was introduced in [15]. Other more recent approaches further adapt the procedure toward time series analysis, e.g., time series segments of various lengths were used in [2], to allow for efficiently handling warping, while an approach that employs temporal modeling was proposed in [1]. Quite recently, a neural formulation of the BoF model was used to perform time series analysis [33], while an extension of this method, that allows for better capturing the temporal dynamics of time series, was introduced in [32].

In contrast with [32], in this work a logistic Neural BoF formulation is used. This allows for training temporal BoF models without using any sophisticated initialization schemes and/or carefully tuning any hyper-parameter, e.g., the initial scaling factor of the kernel function that was employed in [32]. Furthermore, in this work, we studied behavior of the BoF model when combined with deep feature extractors and we appropriately designed an adaptive scaling method that allows for the smooth flow of information in deep BoF-based architectures. To the best of our knowledge, this is the first work in which a deep temporal formulation of the BoF model is used with deep feature extraction layers, after appropriately adapting it to the needs of the specific application, demonstrating that it is indeed possible to learn powerful deep learning models for time series analysis that outperform other competitive state-of-the-art methods.

# 3 Proposed Method

In this Section, the proposed Temporal Logistic Neural Bag-of-Features formulation (abbreviated as "TLo-NBoF" in the rest of this paper) is derived and adapted to the needs of modeling the temporal dynamics of high frequency limit-order book data. Furthermore, as it was already discussed in Section 1, directly employing a BoF formulation, e.g., N-BoF [30], in deep neural networks requires
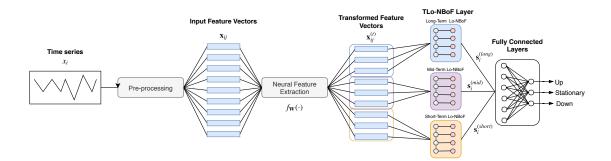
Figure 1: The proposed Temporal Logistic Neural BoF (TLo-NBoF) architecture for time series forecasting.

the careful fine-tuning of several hyper-parameters, e.g., separately adjusting the learning rates per layer, carefully selecting the activation functions and the distribution used for initializing the parameters of the network, etc. In this Section, we delve into these issues, examining some of the reasons for the difficulties that arise when Neural BoF formulations are combined with deep neural networks. Then, the proposed model is appropriately adapted to overcome the aforementioned issues, allowing for directly employing it to learn powerful deep neural network architectures for time series analysis.

## 3.1  Temporal Logistic Neural Bag-of-Features

Let $x_i$ be the $i$-th time series of a collection of $N$ training time series, denoted by $\mathcal{X} = \{x_1, x_2, \ldots, x_N\}$. Then, several feature vectors can be extracted from each time series using several different approaches, as proposed in the literature [16, 21, 1, 10]. Perhaps the most straightforward one is to directly use the raw time series data for each time step as a separate vector [16]. Depending on the application, more sophisticated methods have been also proposed. For example, when dealing with financial data, domain knowledge can be used to design and extract more rich features that describe several aspects of the time series, e.g., multiple feature vectors can be extracted from high frequency limit order book data using the approach proposed in [21]. The $j$-th feature vector extracted from the $x_i$ time series is denoted by $\mathbf{x}_{ij} \in \mathbb{R}^D$, where $D$ is the dimensionality of the extracted feature vectors. Each time series can be then described by the set of the extracted feature vectors, i.e., $x_i = \{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \ldots, \mathbf{x}_{iN_i}\}$, where $N_i$ is the length of the $i$-th time series. Note that different time series might have different lengths, so the proposed method must be able to handle objects that are

5

composed of a varying number of feature vectors.

The extracted features are then transformed using a series of neural transformation layers, as shown in Fig. 1. The employed neural feature extractor is denoted by $f_W(\cdot)$, where $W$ are the parameters (weights) of the feature extraction layers. Any (differential) feature extractor can be used to this end, e.g., convolutional layers [38], recurrent layers [39, 9], etc. In this work, 1-D convolutional layers are used to extract higher level features that are able to capture the temporal relationships between succeeding feature vectors and, as a result, better model the dynamics of a time series. After these neural feature extraction layers, each time series can be represented by a set of *transformed higher level features*. These features are denoted by $\mathbf{x}_{ij}^{(t)} = f_W(x_i, j) \in \mathbb{R}^D$, where $D$ the dimensionality of the transformed features. In this work $D$ equals to the number of filters used in the last convolutional layer, since a convolutional layer is used to transform the input features. Therefore, after the neural feature extraction process, the $i$-th time series is described by $x_i^{(t)} = \{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \ldots, \mathbf{x}_{iN_i}\}$.

Even though the extracted feature vectors $\mathbf{x}_{ij}^{(t)}$ capture higher level information regarding the dynamics of the time series, it is not possible to directly use them for classification purposes (or any other data analysis task), since they have first to be aggregated into a representation that has constant length and it is invariant to the length of the input time series. The performance and flexibility of the resulting model critically depends on the aggregation method that will be employed. For example, the most straightforward approach is to directly flatten the extracted features into one long vector and then feed this vector into a classifer. However, this approach does not allow for a) handling time series that have different lengths, since the length of the resulting vector depends on the length of each series, and b) severely limits the ability of the model to handle time-wrapping/temporal translations/etc.

To overcome the aforementioned limitation, a Temporal Logistic Neural Bag-of-Features formulation is used to efficiently aggregated the extracted feature vectors. Let $\mathcal{V} = \{\mathbf{v}_1, \mathbf{v}_2, \ldots, \mathbf{v}_{N_K}\}$ be a dictionary of $N_K$ codewords that are used to quantized the feature vectors extracted using the neural feature extraction layers. Each codeword is denoted by a vector $\mathbf{v}_k \in \mathbb{R}^D$. Traditionally, these codewords are either selected using the k-means algorithm [18], or by directly sampling them from the set of the extracted feature vectors [3]. However, these approaches cannot be used when the BoF model is combined with deep neural networks, since the input feature distribution is not stationary. Instead, it is constantly shifting requiring the codebook to be updated after each

training step. To this end, the codewords are considered part of the trainable parameters of the network and directly learned using the regular back-propagation algorithm (as it will be demonstrated later). However, the aforementioned processes (clustering or random sampling) can be still used for initializing the codewords [30].

Assuming that each transformed feature vector is generated independently and identically distributed from an unknown distribution controlled by the (image-specific) the vector $\mathbf{s}_i = (s_{i1}, s_{i2}, ..., s_{iN_K})$, then the probability of observing a transformed feature vector $\mathbf{x}_{ij}^{(t)}$ given the $i$-th time series can be estimated using Kernel Density Estimation [34, 20]:

$$p(\mathbf{x}_{ij}^{(t)}|x_i) = \sum_{k=1}^{N_K} s_{ik} K(\mathbf{x}_{ij}^{(t)}, \mathbf{v}_k), \tag{1}$$

where $K(\cdot)$ is a kernel and $s_{ik}$ are the time series specific parameters that control the density estimation. The parameter vector $\mathbf{s}_i$ can be estimated using a maximum likelihood estimator:

$$\mathbf{s}_i = \arg\max_{\mathbf{s}} \sum_{j=1}^{N_i} \log\left(\sum_{k=1}^{N_K} s_{ik} K(\mathbf{x}_{ij}^{(t)}, \mathbf{v}_k)\right). \tag{2}$$

It can be easily derived, as also shown in [3], that these parameters can be effectively estimated as:

$$s_{ik} = \frac{1}{N_i} \sum_{j=1}^{N_i} u_{ijk}, \tag{3}$$

where

$$u_{ijk} = \frac{K(\mathbf{x}_{ij}^{(t)}, \mathbf{v}_k)}{\sum_{l=1}^{N_K} K(\mathbf{x}_{ij}^{(t)}, \mathbf{v}_l)}, \tag{4}$$

giving rise to the well known BoF with soft-assignments [40, 30]. Therefore, the vector $\mathbf{s}_i$ is a histogram that describes the behavior of the $i$-th time series and can be used for the subsequent classification tasks, apart from using it for estimating the distribution of the feature vectors for each time series. Furthermore, as shown in our previous work [30], (4) can be directly implemented as a normalized RBF layer, while (3) can be implemented as a recurrent accumulation layer. Usually a Gaussian kernel is used to calculate the normalized membership vector expressed by (4):

$$K(\mathbf{x}, \mathbf{v}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-||\mathbf{x} - \mathbf{v}||_2^2}{2\sigma^2}\right) \tag{5}$$

7

where $\mathbf{x}$ is a feature vector, $\mathbf{v}$ is a codeword and $\sigma$ is the width of the kernel. However, as shown in [30, 31], it is not straightforward to select the appropriate kernel width and the performance of the model critically relies on the selection of this parameter (even when $\sigma$ is optimized during the training process). To overcome this limitation, in this paper we propose replacing the Gaussian kernel with a more well-behaved and easy to use sigmoid (also known as hyperbolic) kernel [8]:

$$K(\mathbf{x}, \mathbf{v}) = \tanh(\alpha \mathbf{x}^T \mathbf{v} + \beta), \tag{6}$$

where $\alpha$ and $\beta$ are the parameters of the kernel (usually set to $\alpha = 1$ and $\beta = 0$) and $tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$. The kernel is also scaled to $0 \ldots 1$ to ensure that it is compatible with the quantization process employed by (4):

$$K(\mathbf{x}, \mathbf{v}) = \frac{1}{2} \left( \tanh(\alpha \mathbf{x}^T \mathbf{v} + \beta) + 1 \right) = \mathrm{sigm}(2\alpha \mathbf{x}^T \mathbf{v} + 2\beta), \tag{7}$$

where $sigm(x) = \frac{1}{1+e^{-x}}$ is the logistic sigmoid function. Using this kernel also allows for avoiding the need for sophisticated initialization schemes based on computationally intensive algorithms, e.g., k-means, allowing for simply randomly initializing the codebook, along with the other parameters of the network.
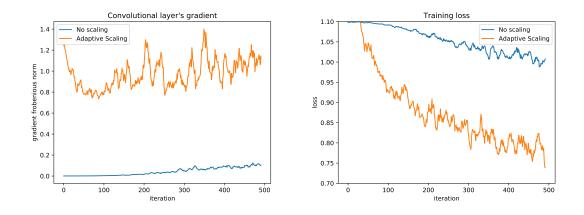


Figure 2: Effect of using *adaptive scaling*, i.e., allowing the network to adjust the scale of the $\mathbf{s}_i$ and $\mathbf{x}_{ij}$ vectors during the training process, on the gradients of the layers before the proposed TLo-NBoF layer (left figure) and on the loss function (right figure) during the first 500 training iterations. The proposed adaptive scaling approach significantly improves the speed of the convergence of the network.

8

Note that the histogram vector $\mathbf{s}_i$ captures only the overall behavior of the $i$-th time series. In this work, a fine-grained temporal segmentation scheme is also proposed to capture the temporal dynamics of the time series. To this end, the transformed feature vectors are segmented into $N_T$ temporal regions, as shown in Fig. 1, to capture the short-term, mid-term, and long-term behavior of the time series ($N_T = 3$ temporal regions are used). Therefore, the most recent $\lfloor \frac{N_i}{N_T} \rfloor$ feature vectors are employed for calculating the short-term histogram $\mathbf{s}_i^{(short)}$, the preceding $\lfloor \frac{N_i}{N_T} \rfloor$ are used to calculate the mid-term histogram $\mathbf{s}_i^{(mid)}$, while the rest of the feature vectors are used for calculating the long-term histogram $\mathbf{s}_i^{(long)}$. Note that different temporal segmentation schemes can be also applied, i.e., the short-term feature vectors can be fed to both the mid-term and long-term Lo-NBoF models to model the corresponding behavior over longer periods of time. Finally, the resulting concatenated vector $\mathbf{s}_i = [\mathbf{s}_i^{(short)}, \mathbf{s}_i^{(mid)}, \mathbf{s}_i^{(long)}] \in \mathbb{R}^{3N_K}$ is fed to the following fully connected layer, as shown in Fig. 1.

## 3.2 Learning Deep Architectures with Temporal Logistic Neural Bag-of-Features

Even though the previously described architecture can deal with time series of variable length and capture their fine-grained temporal dynamics, it requires a significant effort in order to tune the appropriate hyper-parameters, e.g., learning rate, initialization, etc., in order to effectively train the resulting architecture. We argue that the main reason for that is the normalizations involved in (3) and (4). These normalizations, that scales down the $l^1$ norm of the corresponding vectors, prohibit the smooth flow of information, both in the forward and backward propagation. This is illustrated in Fig. 2, where the Frobenius norm of the gradient of the parameters of the first layer of the network are plotted for the first 500 training iterations. Note the extremely small values for the gradients (blue line), that effectively prohibit the gradients from backpropagating to the layers behind the TLo-NBoF model. More than 200 iterations are needed just for starting to slowly update these layers. The harsh scaling involved in (3) and (4) also reduce the variance of the activation/gradients. However, it is well established that maintaining the same variance for these quantities across the various layers of the network is critical to ensure that the network will be correctly trained, as discussed in detail in [11, 12]. Therefore, to overcome these issues, we propose

9

appropriately scaling the $\mathbf{s}_i$ and $\mathbf{u}_{ij}$ vectors as:

$$s_{ik} = c_s \frac{1}{N_i} \sum_{j=1}^{N_i} u_{ijk}, \tag{8}$$

and

$$u_{ijk} = c_u \frac{K(\mathbf{x}_{ij}^{(t)}, \mathbf{v}_k)}{\sum_{l=1}^{N_K} K(\mathbf{x}_{ij}^{(t)}, \mathbf{v}_l)}, \tag{9}$$

where $c_s$ is initialized to $N_K$, while $c_s$ is initialized to the average number of feature vectors per object. Then, the appropriate values for these two scaling factors are learned during the training process, allowing for easily adjusting the norm of the corresponding vectors to better facilitate the training process. Note that a similar approach have been also used by some activation functions, e.g., PReLU [12], to allow the network to better adjust to the task at hand. Note that the scaling that is involved in (8) and (9) still leads to maintaining a constant $l^1$ norm for these vectors, since $c_s$ and $c_u$ are fixed for all the time series that are presented to the network (after the training). This approach, i.e., allowing the network to automatically adjust the norm of the corresponding vectors to allow the smooth flow of information, is called "*adaptive scaling*" through this paper. As it can be shown in Fig. 2, where the training process is illustrated during the first 500 training iterations, adaptive scaling can significantly improve the convergence of the network and allows for have more well-behaved gradients on the layers of the network before the employed TLo-NBoF layer.

The resulting architecture, as shown in Fig 1, can be now directly trained in an end-to-end fashion using gradient descent, i.e.,

$$\Delta(\mathbf{W}_{conv}, \mathbf{V}, \mathbf{W}_{fc}, \mathbf{c}) = \eta(\frac{\partial \mathcal{L}}{\partial \mathbf{W}}, \frac{\partial \mathcal{L}}{\partial \mathbf{V}}, \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{fc}}, \frac{\partial \mathcal{L}}{\partial \mathbf{c}}), \tag{10}$$

where $\mathcal{L}$ is the employed loss function, $\mathbf{W}$ denotes the parameters of the neural feature extractor $f_W(\cdot)$, $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \ldots, \mathbf{v}_{N_K}]$ denotes the codebook used by the proposed model, $\mathbf{W}_{fc}$ denotes the parameters of the fully connected layers and $\mathbf{c} = (c_u, c_s)$ are the scaling parameters involved in adaptive scaling. The cross-entropy loss is used for all the experiments conducted in this paper, while the same codebook is utilized for all the temporal regions. The Adam algorithm is used to perform the optimization [22]. The training time series were fed to the network in batches of 128 samples, where each time series was sampled with probability inversely proportional to the

Table 1: Architecture of the model employed for financial time series forecasting

| Layer | Output size |
|---|---|
| Input | $N_i \times 144$ |
| Convolutional (256 filters, kernel size 5) | $N_i \times 256$ |
| TLo-NBoF ($N_K = 256, N_T = 3$) | $3 \times 256$ |
| Fully Connected (512) | 512 |
| Fully Connected (3) | 3 |

frequency of its class. The learning rate was set to $\eta = 10^{-4}$, while the networks were trained for 20 epochs. Finally, note that the parameters of the kernel $\alpha$ and $\beta$ can be also optimized during the training process:

$$\Delta(\alpha, \beta) = \eta(\frac{\partial \mathcal{L}}{\partial \alpha}, \frac{\partial \mathcal{L}}{\partial \beta}), \tag{11}$$

We refer to this approach, i.e., learning the kernel parameters $\alpha$ and $\beta$, as *"Kernel Parameter Learning"*, in the rest of this paper.

Table 2: Ablation Study (the prediction horizon is set to the next 10 time steps)

| Deep Features | Temp. Model. | Kernel Param. | Adaptive Scaling | Macro-F1 | Cohen's $\kappa$ |
|---|---|---|---|---|---|
| - | - | - | ✓* | $42.66 \pm 0.28$ | $0.1847 \pm 0.0026$ |
| ✓ | - | - | ✓* | $46.77 \pm 1.53$ | $0.2219 \pm 0.0230$ |
| - | ✓ | - | ✓* | $50.14 \pm 1.36$ | $0.2686 \pm 0.0179$ |
| ✓ | ✓ | - | ✓* | $51.65 \pm 0.99$ | $0.2783 \pm 0.0109$ |
| ✓ | ✓ | ✓ | - | $50.65 \pm 0.71$ | $0.2603 \pm 0.0119$ |
| ✓ | ✓ | ✓ | ✓* | $53.48 \pm 0.45$ | $0.3013 \pm 0.0075$ |
| ✓ | ✓ | ✓ | ✓ | $\mathbf{53.54 \pm 0.24}$ | $\mathbf{0.3031 \pm 0.0066}$ |

(✓* refers to using the scaling parameters $c_s = N_K$ and $c_u = E[N_i]$, but not adjusting them during the training process)

# 4  Experimental Evaluation

The proposed method was evaluated using a large-scale limit order book dataset [27, 26]. The employed dataset consists of high frequency limit order book data collected from 5 Finish companies traded in the Helsinki Exchange (operated by Nasdaq Nordic). The 10 highest and lower ask order prices were collected for each time step, while data were collected over a period of 10 business days

Table 3: Evaluation results using the FI-2010 dataset

| Method | Macro Precision | Macro Recall | Macro F1 score | Cohen's $\kappa$ |
|---|---|---|---|---|
| MLP [32] | $40.20 \pm 0.50$ | $56.25 \pm 2.20$ | $36.91 \pm 1.81$ | $0.1281 \pm 0.0137$ |
| BoF [32] | $39.26 \pm 0.94$ | $51.44 \pm 2.53$ | $36.28 \pm 2.85$ | $0.1182 \pm 0.0246$ |
| N-BoF [32] | $42.28 \pm 0.87$ | $61.41 \pm 3.68$ | $41.63 \pm 1.90$ | $0.1724 \pm 0.0212$ |
| T-BoF [32] | $43.85 \pm 1.11$ | $66.66 \pm 3.40$ | $43.96 \pm 1.59$ | $0.1992 \pm 0.0201$ |
| WMTR [37] | $46.25 \pm N/A$ | $51.29 \pm N/A$ | $47.87 \pm N/A$ | $N/A$ |
| CNN (256 filters) | $44.69 \pm 1.13$ | $58.70 \pm 1.85$ | $47.19 \pm 1.69$ | $0.2192 \pm 0.0235$ |
| LSTM (256 neurons) | $47.63 \pm 2.25$ | $52.60 \pm 2.74$ | $49.51 \pm 2.43$ | $0.2395 \pm 0.0388$ |
| GRU (256 neurons) | $47.70 \pm 2.09$ | $56.76 \pm 2.99$ | $50.55 \pm 2.33$ | $0.2560 \pm 0.0364$ |
| TLo-NBoF (proposed) | $\mathbf{50.20 \pm 2.22}$ | $\mathbf{58.19 \pm 2.13}$ | $\mathbf{52.98 \pm 2.37}$ | $\mathbf{0.2900 \pm 0.0361}$ |

(1st June 2010 to 14th June 2010). A total of 4.5 million limit orders were gathered and processed according to the pre-processing and feature extraction pipeline proposed in [21]. Thus, a total number of 453,975 144-dimensional feature vectors were extracted.

The anchored evaluation setup proposed in [36] was used for the evaluation: The time series that were extracted from the first day were used to train the model, while the data from the second day were employed for the evaluation. Then, the first two days were used for the training and the next day was used for the evaluation, etc. This process was repeated 9 times. For all the evaluated metrics (macro-precision, macro-recall, macro-F1 and Cohen's $\kappa$ [6]), the mean and standard deviation are reported. The direction of the average mid price (up, stationary or down) after 10 time steps was predicted. A stock was considered to be stationary if the change in the mid price was less than to 0.01%.

For each time step a time series that consists of the last 15 feature vectors was compiled. The time series was segmented into $N_T = 3$ temporal regions, each consisting of 5 feature vectors. The detailed architecture of the employed network is shown in Table 1. First, a convolutional feature extractor with 256 filters (kernel size was set to 5) was used. Then, the transformed feature vectors were fed to a TLo-NBoF layer with 256 codewords. Finally, the temporal histograms extracted from the TLo-NBoF layer were fed into two fully connected layers responsible for predicting the future behavior of the price for the given time series. The ReLU function was used both for the convolutional feature extractor and the first fully connected layer. Note that the TLo-NBoF layer can be directly implemented by using 1D convolutions with kernel size 1, setting the weights of the convolution equal to the codebook and then using the appropriate activation and scaling layers.

First, the effect of the different parts of the proposed architecture are evaluated in the ablation

study provided in Table 2. The last three days of the training data were used for performing the ablation study. The first line refers to using a plain Lo-NBoF layer with the appropriate scaling ($c_s$ and $c_u$) to ensure that the model will be successfully trained. Then, adding a convolutional feature extractor ("Deep Features") improves the Cohen's $\kappa$ by 20%, while using temporal modeling, i.e., three separate histograms that describe the short-term, mid-term and long-term behavior, improves the Cohen's $\kappa$ by over 40%. Combining the deep feature extractor with the temporal modeling further improves the performance of the proposed model. Learning the parameters of the kernel (Kernel Parameter Learning) further boost the metrics. Finally note that using and learning the scaling parameters $c_s$ and $c_u$ is crucial for successfully training the network, since without them the performance of the models is reduced, e.g., Cohen's $\kappa$ is reduced by over 14%.

The proposed method is also compared to various other baselines proposed in the literature [37, 32], as well as other powerful convolutional and recurrent deep learning models. For the CNN baseline the same architecture as the one shown in Table 1 was used, but the TLo-NBoF layer was replaced by a Global Average Pooling layer. For the GRU [19] and LSTM [14] models, the feature extraction layers and TLo-NBoF layer were replaced by the appropriate recurrent model. The final state of these models was used for performing the classification. The proposed method significantly improves the performance metrics over both the plain Temporal BoF model [32] and the more powerful recurrent and convolutional architectures (the performance is improved by over 13% over the next best performing model according to the $\kappa$ metric).

# 5    Conclusions

In this paper, a novel logistic formulation of the Neural Bag-of-Features model was proposed and appropriately adapted in order to be efficiently used with deep feature extractors. In this way, the proposed method can effectively combine the advantages of the BoF model with the great learning capacity of DL models, leading to powerful models for time series analysis. The employed fully differential logistic formulation of the BoF model, together with the proposed adaptive scaling mechanism, allows for directly training the resulting architecture in an end-to-end fashion. Furthermore, the proposed method is capable of modeling the behavior of time series at various temporal levels. Finally, the proposed method was extensively evaluated using a large-scale financial time series dataset that consists of more than 4 million limit orders and it was demonstrated that it

performs better than other competitive baseline and state-of-the-art methods.

# References

[1] Adeline Bailly, Simon Malinowski, Romain Tavenard, Thomas Guyet, and Laetitia Chapel. Bag-of-temporal-sift-words for time series classification. In *ECML/PKDD Workshop on Advanced Analytics and Learning on Temporal Data*, 2015.

[2] Mustafa Gokce Baydogan, George Runger, and Eugene Tuv. A bag-of-features framework to classify time series. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(11):2796–2802, 2013.

[3] Subhabrata Bhattacharya, Rahul Sukthankar, Rong Jin, and Mubarak Shah. A probabilistic representation for efficient large scale visual recognition tasks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2593–2600, 2011.

[4] Émilie Poisson Caillault, Alain Lefebvre, André Bigand, et al. Dynamic time warping-based imputation for univariate time series data. *Pattern Recognition Letters*, 2017.

[5] Li-Juan Cao and Francis Eng Hock Tay. Support vector machine with adaptive parameters in financial time series forecasting. *IEEE Transactions on Neural Networks*, 14(6):1506–1518, 2003.

[6] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46, 1960.

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

[8] Inderjit S Dhillon, Yuqiang Guan, and Brian Kulis. Kernel k-means: spectral clustering and normalized cuts. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 551–556, 2004.

[9] Matthew Dixon. Sequence classification of the limit order book using recurrent neural networks. *Journal of Computational Science*, 24:277 – 286, 2018.

[10] Zhouyu Fu, Guojun Lu, Kai Ming Ting, and Dengsheng Zhang. Music classification via the bag-of-features approach. *Pattern Recognition Letters*, 32(14):1768–1777, 2011.

[11] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, pages 249–256, 2010.

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1026–1034, 2015.

[13] Henrique Steinherz Hippert, Carlos Eduardo Pedreira, and Reinaldo Castro Souza. Neural networks for short-term load forecasting: A review and evaluation. *IEEE Transactions on Power Systems*, 16(1):44–55, 2001.

[14] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

[15] Alexandros Iosifidis, Anastasios Tefas, and Ioannis Pitas. Dynamic action recognition based on dynemes and extreme learning machine. *Pattern Recognition Letters*, 34(15):1890–1898, 2013.

[16] Alexandros Iosifidis, Anastasios Tefas, and Ioannis Pitas. Multidimensional sequence classification based on fuzzy distances and discriminant analysis. *IEEE Transactions on Knowledge and Data Engineering*, 25(11):2564–2575, 2013.

[17] Alexandros Iosifidis, Anastastios Tefas, and Ioannis Pitas. Discriminant bag of words based representation for human action recognition. *Pattern Recognition Letters*, 49:185–192, 2014.

[18] Anil K Jain. Data clustering: 50 years beyond k-means. *Pattern recognition letters*, 31(8):651–666, 2010.

[19] Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. An empirical exploration of recurrent network architectures. In *Proceedings of the International Conference on Machine Learning*, pages 2342–2350, 2015.

[20] Vladimir Katkovnik and Ilya Shmulevich. Kernel density estimation with adaptive varying window size. *Pattern Recognition Letters*, 23(14):1641–1648, 2002.

[21] Alec N. Kercheval and Yuan Zhang. Modelling high-frequency limit order book dynamics with support vector machines. *Quantitative Finance*, 15(8):1315–1329, 2015.

[22] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations*, 2015.

[23] Martin Längkvist, Lars Karlsson, and Amy Loutfi. A review of unsupervised feature learning and deep learning for time-series modeling. *Pattern Recognition Letters*, 42:11–24, 2014.

[24] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436, 2015.

[25] Zachary C Lipton, David C Kale, Charles Elkan, and Randall Wetzell. Learning to diagnose with lstm recurrent neural networks. *arXiv preprint 1511.03677*, 2015.

[26] Paraskevi Nousi, Avraam Tsantekidis, Nikolaos Passalis, Adamantios Ntakaris, Juho Kanniainen, Anastasios Tefas, Moncef Gabbouj, and Alexandros Iosifidis. Machine learning for forecasting mid price movement using limit order book data. *arXiv preprint arXiv:1809.07861*.

[27] Adamantios Ntakaris, Martin Magris, Juho Kanniainen, Moncef Gabbouj, and Alexandros Iosifidis. Benchmark dataset for mid-price prediction of limit order book data. *arXiv preprint arXiv:1705.03233*, 2017.

[28] Nikolaos Passalis and Anastasios Tefas. Entropy optimized feature-based bag-of-words representation for information retrieval. *IEEE Transactions on Knowledge and Data Engineering*, 28(7):1664–1677, 2016.

[29] Nikolaos Passalis and Anastasios Tefas. Learning bag-of-features pooling for deep convolutional neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017.

[30] Nikolaos Passalis and Anastasios Tefas. Neural bag-of-features learning. *Pattern Recognition*, 64:277–294, 2017.

[31] Nikolaos Passalis and Anastasios Tefas. Training lightweight deep convolutional neural networks using bag-of-features pooling. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–11, 2018.

[32] Nikolaos Passalis, Anastasios Tefas, Juho Kanniainen, Moncef Gabbouj, and Alexandros Iosifidis. Temporal bag-of-features learning for predicting mid price movements using high frequency limit order book data. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2018.

[33] Nikolaos Passalis, Avraam Tsantekidis, Anastasios Tefas, Juho Kanniainen, Moncef Gabbouj, and Alexandros Iosifidis. Time-series classification using neural bag-of-features. In *Proceedings of the European Signal Processing Conference*, pages 301–305, 2017.

[34] Bernard W Silverman. *Density estimation for statistics and data analysis*. Routledge, 2018.

[35] Josef Sivic, Andrew Zisserman, et al. Video google: A text retrieval approach to object matching in videos. In *Proceedings of the IEEE International Conference on Computer Vision*, volume 2, pages 1470–1477, 2003.

[36] Emilio Tomasini and Urban Jaekle. *Trading Systems*. Harriman House Limited, Hampshire, UK, 2011.

[37] Dat Thanh Tran, Martin Magris, Juho Kanniainen, Moncef Gabbouj, and Alexandros Iosifidis. Tensor representation in high-frequency financial data for price change prediction. In *Proceedings of the IEEE Symposium Series on Computational Intelligence*, pages 1–7, 2017.

[38] Avraam Tsantekidis, Nikolaos Passalis, Anastasios Tefas, Juho Kanniainen, Moncef Gabbouj, and Alexandros Iosifidis. Forecasting stock prices from the limit order book using convolutional neural networks. In *Proceedings of the IEEE Conference on Business Informatics (CBI)*, pages 7–12, 2017.

[39] Avraam Tsantekidis, Nikolaos Passalis, Anastasios Tefas, Juho Kanniainen, Moncef Gabbouj, and Alexandros Iosifidis. Using deep learning to detect price change indications in financial markets. In *Proceedings of the European Signal Processing Conference*, pages 2511–2515, 2017.

[40] Jan C Van Gemert, Jan-Mark Geusebroek, Cor J Veenman, and Arnold WM Smeulders. Kernel codebooks for scene categorization. In *Proceedings of the European Conference on Computer Vision*, pages 696–709, 2008.

[41] Yi Wang, Zhiming Luo, and Pierre-Marc Jodoin. Interactive deep learning method for segmenting moving objects. *Pattern Recognition Letters*, 96:66–75, 2017.