# Attention-based Neural Bag-of-Features Learning for Sequence Data

Dat Thanh Tran*, Nikolaos Passalis†, Anastasios Tefas†, Moncef Gabbouj*, Alexandros Iosifidis‡

*Department of Computing Sciences, Tampere University, Tampere, Finland
†Department of Informatics, Aristotle University of Thessaloniki, Thessaloniki, Greece
‡Department of Engineering, Aarhus University, Aarhus, Denmark
thanh.tran@tuni.fi, passalis@csd.auth.gr, tefas@aiia.csd.auth.gr,
moncef.gabbouj@tuni.fi, alexandros.iosifidis@eng.au.dk

*Abstract*—In this paper, we propose 2D-Attention (2DA), a generic attention formulation for sequence data, which acts as a complementary computation block that can detect and focus on relevant sources of information for the given learning objective. The proposed attention module is incorporated into the recently proposed Neural Bag of Feature (NBoF) model to enhance its learning capacity. Since 2DA acts as a plug-in layer, injecting it into different computation stages of the NBoF model results in different 2DA-NBoF architectures, each of which possesses a unique interpretation. We conducted extensive experiments in financial forecasting, audio analysis as well as medical diagnosis problems to benchmark the proposed formulations in comparison with existing methods, including the widely used Gated Recurrent Units. Our empirical analysis shows that the proposed attention formulations can not only improve performances of NBoF models but also make them resilient to noisy data.

## I. Introduction

Learning problems in many fields involve sequence data such as time-series forecasting [1], [2], audio analysis [3], [4] or natural language processing [5], [6], all of which have been extensively studied. In many application scenarios, the observed sequence is highly non-stationary and noisy, which makes the task of modeling the underlying generating process more difficult. For example, in sound source separation in which the objective is to recover different unknown sources by filtering the observed mixtures, the existence of environmental noise is inherent and often complicates the separation process. Several mathematical techniques have been proposed to model the underlying data and noise distributions or to extract hand-crafted features, capturing certain desirable properties. In financial time-series analysis, representative examples include autoregressive (AR) and moving average (MA) [7] features, which were later extended with a differencing step to eliminate nonstationarity, known as autoregressive integrated moving average (ARIMA) [8]. Gaussian processes and Hidden Markov Model were popular mathematical frameworks in audio analysis. To ensure mathematical and computational tractability, these classical models are often formulated under many assumptions, which are sensitive to initialization and misaligned with real-world conditions, thus limiting their professional usage in practice.

During the last decade, thanks to the development in stochastic optimization techniques and computing hardware, as well as the declining costs of data acquisition and storage, a data-driven approach based on deep neural networks and stochastic optimization has replaced the classical model-based approach and convex optimization. Nowadays, many of the state-of-the-art solutions for learning with sequence data are developed on the basis of neural networks. Notably, a class of neural network architecture called Recurrent Neural Networks (RNN), which is specifically designed to process variable-length sequences and to capture sequential patterns, has become the main workforce in different application domains. Another dedicated neural formulation for sequence data is the bilinear structures [1], [9], [10], which were proposed to separately capture the dependencies along the temporal and spatial dimension in financial time-series. Even existing neural architectures, which were originally proposed for visual inputs such as Convolutional Neural Network (CNN) [11] and Neural Bag-of-Features (NBoF) [12], have shown competitive performances in tackling sequence data compared to dedicated statistical models [13], [14]. The advantage of neural formulations over statistical learning and traditional hand-crafted features lies in the fact that fewer assumptions are made, and data is leveraged to automatically identify and extract task-relevant features in an end-to-end fashion.

Bag-of-Features (BoF) model [15] was originally proposed to build histogram representations from images. Later, it was shown that BoF could be successfully applied to extract high-level representations for other data modalities such as video and audio [16], [17], [18], [19]. Learning BoF representations consists of two steps: *dictionary learning* and *feature quantization and encoding*. In the dictionary learning step, each object is first represented by a set of low-level features, which could be, for example, a collection of local descriptors like SIFT [20] for image object or word-level vector-encodings for a sentence object. These features are then used to generate a compact dictionary (codebook) comprising of the most representative features, also known as *codewords*. In the second step, the histogram representation of each object is extracted by quantizing its low-level features using the codebook.

Recently, Neural Bag-of-Features (NBoF) [12], a neural network generalization of the BoF model, has been proposed. Similar to its predecessor, NBoF can generate a fixed-size histogram vector from variable-size inputs. This neural network generalization works as a feature extraction layer, which can be combined and optimized jointly with other neural network layers to tackle both unsupervised and supervised objectives via stochastic optimization. Since the dictionary learning step in NBoF is updated in conjunction with other layers towards the end goal of optimizing an objective function, histogram vectors synthesized by NBoF are more representative than those produced by BoF in different learning scenarios such as visual recognition, information retrieval, and financial forecasting [21], [12], [14].

While the NBoF model works well in different learning problems, the current formulations still possess some limitations. In the aggregation step, all of the quantized features are simply averaged to form the histogram vector. For sequence data, this implies that the model only allows equal contributions of the quantized features coming from different time steps to form the output representation. Similarly, the quantization results produced by each codeword are considered equally important for every sequence in the training set. These properties limit the dictionary learning, quantization, and encoding process to fully take advantage of the data-driven approach.

To incorporate a higher degree of flexibility into the NBoF model, a weighing mechanism on a sequence level is desirable. That is, for each individual sequence, the model has the flexibility to perform a weighted sum of quantized outputs in the aggregation step, with the coefficients being adaptively changed with respect to the input sequence, or to select/discard irrelevant codewords, given the input sequence. In neural network literature, this is often achieved by having some attention mechanisms [22], [23], [10]. The idea of attention is inspired by the phenomenon observed in the human visual cortex that visual stimuli from multiple objects actively compete for neural encoding.

Although various attention mechanisms have been proposed for existing neural network architectures such as CNN [22], [24], LSTM [23], [25] or Bilinear structure [10], there is yet any formulation for the NBoF model when learning with sequence data. To have a generic attention mechanism that can be applied in a plug-and-play manner, in this work, we propose 2D-Attention (2DA), a neural network module that promotes competitions among different rows or columns in the input matrix and only (soft) selects those which win for attention. We will then demonstrate that by injecting 2DA into NBoF, we can overcome those limitations mentioned previously. The contributions of our work can be summarized as follows:

- We propose a new type of attention formulation for matrix data, which is dubbed as 2DA. The proposed layer acts as a complementary computation block, which is capable of identifying relevant sources of information to perform selective masking on the given input matrix.

- We incorporate 2DA into different stages of the Neural Bag-of-Features (NBoF) model, creating various 2DA-NBoF extensions that can enhance the feature quantization or histogram accumulation step in the NBoF model. Extensive experiments were conducted in three different application domains: financial forecasting, audio analysis, and medical diagnosis, which demonstrate the effectiveness of our attention module in improving the NBoF model. In cases of noisy input, a variant of 2DA-NBoF shows resilience to noises by filtering out the noisy source of information before the feature quantization step.

The remainder of the paper is organized as follows: in Section II, we review the NBoF model and its extensions for time-series data, as well as previously proposed attention mechanisms in the neural network literature. In Section III, we first present the proposed attention module 2DA and its interpretation. Several extensions of the NBoF model that incorporates 2DA are then presented. In Section IV, we provide details of our experiment protocols and quantitative analysis. Section V concludes our work with possible future research directions.

## II. RELATED WORK

The NBoF model [12] consists of two components: a quantization layer and an accumulation layer. Each quantization neuron in the quantization layer performs like a codeword, which can be updated via BackPropagation algorithm. In the original formulation [12], the Radial Basis Function (RBF) layer was used for feature quantization. Recently, it has been shown that the hyperbolic kernel is also effective for the feature quantization step [26]. Here we describe the original formulation with RBF layer.

Let $K$ be the number of neurons (codewords) in the RBF layer and $\mathbf{v}_k \in \mathbb{R}^D$ be the $k$-th codeword. In addition, the shape of the Gaussian function modeled by each neuron can be adjusted via parameter $\mathbf{w}_k$. Let us denote the sequence of $N$ features as $\mathbf{X} = [\mathbf{x}_1, \ldots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ with $\mathbf{x}_n \in \mathbb{R}^D, n = 1, \ldots, N$. The output of the $k$-th RBF neuron given the input feature $\mathbf{x}_n$ is the following:

$$\phi_{n,k} = \frac{\exp\big(-\|(\mathbf{x}_n - \mathbf{v_k}) \odot \mathbf{w}_k\|_2\big)}{\sum_{m=1}^{K} \exp\big(-\|(\mathbf{x}_n - \mathbf{v}_m) \odot \mathbf{w}_m\|_2\big)} \quad (1)$$

where $\odot$ denotes element-wise product and $\mathbf{w}_k \in \mathbb{R}^D$ is the learnable weight vector that enables the shape of Gaussian kernel associated with the $k$-th RBF neuron to change.

As the sequence $\mathbf{X}$ goes through the quantization layer, each feature $\mathbf{x}_n$ is quantized as $\boldsymbol{\phi}_n = [\phi_{n,1}, \ldots, \phi_{n,K}]^T \in \mathbb{R}^K$, producing a sequence of quantized features $\boldsymbol{\Phi} = [\boldsymbol{\phi}_1, \ldots, \boldsymbol{\phi}_N] \in \mathbb{R}^{K \times N}$. The accumulation layer aggregates the information in $\boldsymbol{\Phi}$ by calculating the averaged quantized feature:

$$\mathbf{y} = \frac{1}{N} \sum_{n=1}^{N} \boldsymbol{\phi}_n \quad (2)$$

There have been few extensions of NBoF model for sequence data. For example, Temporal Neural Bag-of-Features

(TNBoF) model with different specialized codebooks has been proposed in [27] to capture both short-term and long-term temporal information in financial time-series. In [26], the authors derived the logistic formulation of the NBoF model using the hyperbolic kernel instead of the RBF kernel for the quantization step and proposed an adaptive scaling mechanism which showed significant improvements in training stability and performance of the NBoF networks.

While the attention mechanism was biologically inspired from the perspective of visual processing, this technique has also inspired and advanced several works in sequence data analysis, notably in sequence-to-sequence learning tasks. The first attention formulation applied to sequence data was proposed in [6] for tackling machine translation tasks. In this formulation, the authors proposed to construct the context vectors in Sequence-to-sequence Recurrent Neural Network model by selectively combining some hidden states, rather than using the last hidden state as the context vector. The selection coefficients, also known as attention weights, are computed adaptively based on the given input sequence, and updated jointly with other parameters during stochastic optimization.

The successful application of attention mechanism in machine translation tasks has led to the emergence of other attention formulations, which are designed to capture different types of salient information in sequence data. For example, in [28], the authors proposed a formulation that can detect pseudo-periods in certain types of time-series, such as energy consumption or meteorology data. To predict the future stock index, a dual-stage attention mechanism was proposed in [25] for RNN to actively select relevant exogenous series and temporal instances. Similarly, to highlight and focus on important temporal events in Limit Order Book, the authors in [10] proposed a method to calculate attention masks for bilinear networks. Although an attention formulation has been proposed for the convolutional NBoF model in [24] to estimate the true color of images capturing by different devices, this formulation only works with image data. To the best of our knowledge, there has been no attention formulation for the NBoF model to tackle sequence data.

## III. PROPOSED METHODS

In this Section, we will first present 2D-Attention (2DA), our proposed attention calculation for matrix data. Then, we will show how 2DA can be used to address different limitations of the NBoF model as described in Section I. Throughout the paper, we denote scalar values by either lower-case or upper-case characters $(a, b, A, B, \dots)$, vectors by lower-case bold-face characters $(\mathbf{x}, \mathbf{y}, \dots)$, matrices by upper-case bold-face characters $(\mathbf{X}, \mathbf{Y}, \dots)$, and mathematical functions by calligraphy characters $\mathcal{F}, \mathcal{G}, \dots$. In addition, we use $x_{mn}$ to denote the element at position $(m, n)$ in a matrix $\mathbf{X}$.

### A. 2D-Attention

A matrix $\mathbf{S} \in \mathbb{R}^{M \times N}$ is a second-order tensor which has two modes, with $M$ and $N$ are the dimensions of the first and second mode, respectively. The matrix representation provides a natural way to represent a signal with two different sources of information. For example, a multivariate time-series is represented by a matrix with one mode representing the temporal dimension, and the other mode represents different sources that generate individual series.

The general idea of attention mechanism is to highlight important elements in the data while discarding irrelevant ones. For data represented as a matrix $\mathbf{S}$, rather than considering each element in $\mathbf{S}$ individually, we would like to actively select certain columns or rows of $\mathbf{S}$ while discarding the others. This is because columns or rows of $\mathbf{S}$ usually form coherent sub-groups of the data. For example, discarding some temporal events or some individual series in a multivariate series corresponds to removing some rows or columns, depending on the orientation of the matrix.

To adaptively determine and focus on different columns or rows of a matrix, we propose 2D-Attention (2DA), with the functional form denoted by $\mathcal{F}_{2DA}$. This function takes a matrix $\mathbf{S} \in \mathbb{R}^{M \times N}$ as the input, and returns $\tilde{\mathbf{S}} \in \mathbb{R}^{M \times N}$ as the output. That is:

$$\tilde{\mathbf{S}} = \mathcal{F}_{2DA}(\mathbf{S}) \tag{3}$$

$\tilde{\mathbf{S}}$ can be considered as a filtered version of $\mathbf{S}$, where irrelevant columns of $\mathbf{S}$ with respect to the learning problem are zeroed out. Here we should note that $\mathcal{F}_{2DA}$ performs adaptive attention with respect to columns of $\mathbf{S}$. To focus on different rows of $\mathbf{S}$, we can simply apply $\mathcal{F}_{2DA}$ to the transpose of $\mathbf{S}$.

The selection or rejection of the columns of $\mathbf{S}$ is conducted via element-wise matrix multiplications as follows:

$$\tilde{\mathbf{S}} = \tau(\mathbf{S} \odot \mathbf{A}) + (1 - \tau)\mathbf{S} \tag{4}$$

where $\mathbf{A} \in \mathbb{R}^{M \times N}$ denotes the attention mask with values in the range $[0, 1]$. Each column in $\mathbf{A}$ encodes the importance of the corresponding column in $\mathbf{S}$. That is, the attention mask contains values that are close to $1$ corresponding to those columns in $\mathbf{S}$ that contain important information for the downstream learning task and vice versa. In Eq.(4), parameter $\tau \in \mathbb{R}$, which is jointly optimized with other parameters, is used to allow flexible control of the attention mechanism: when $\mathbf{S}$ contains redundant or noisy information in its columns, the effect of attention mask $\mathbf{A}$ is enabled by pushing $\tau$ close to $1$; on the other hand, when every column of $\mathbf{S}$ is necessary, i.e., there is no need for attention, pushing $\tau$ close to $0$ will disable the effect of $\mathbf{A}$. The necessity of attention is thus automatically determined by optimizing $\tau$ with respect to a given problem.

To calculate the attention mask $\mathbf{A}$, the proposed 2DA method learns to measure the relative importance between columns of $\mathbf{S}$ via a specially designed weight matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$: all elements of $\mathbf{W}$ are learnable, i.e., they are updated during stochastic optimization, except the diagonal elements, which are fixed to $1/N$. The attention mask is calculated as follows:
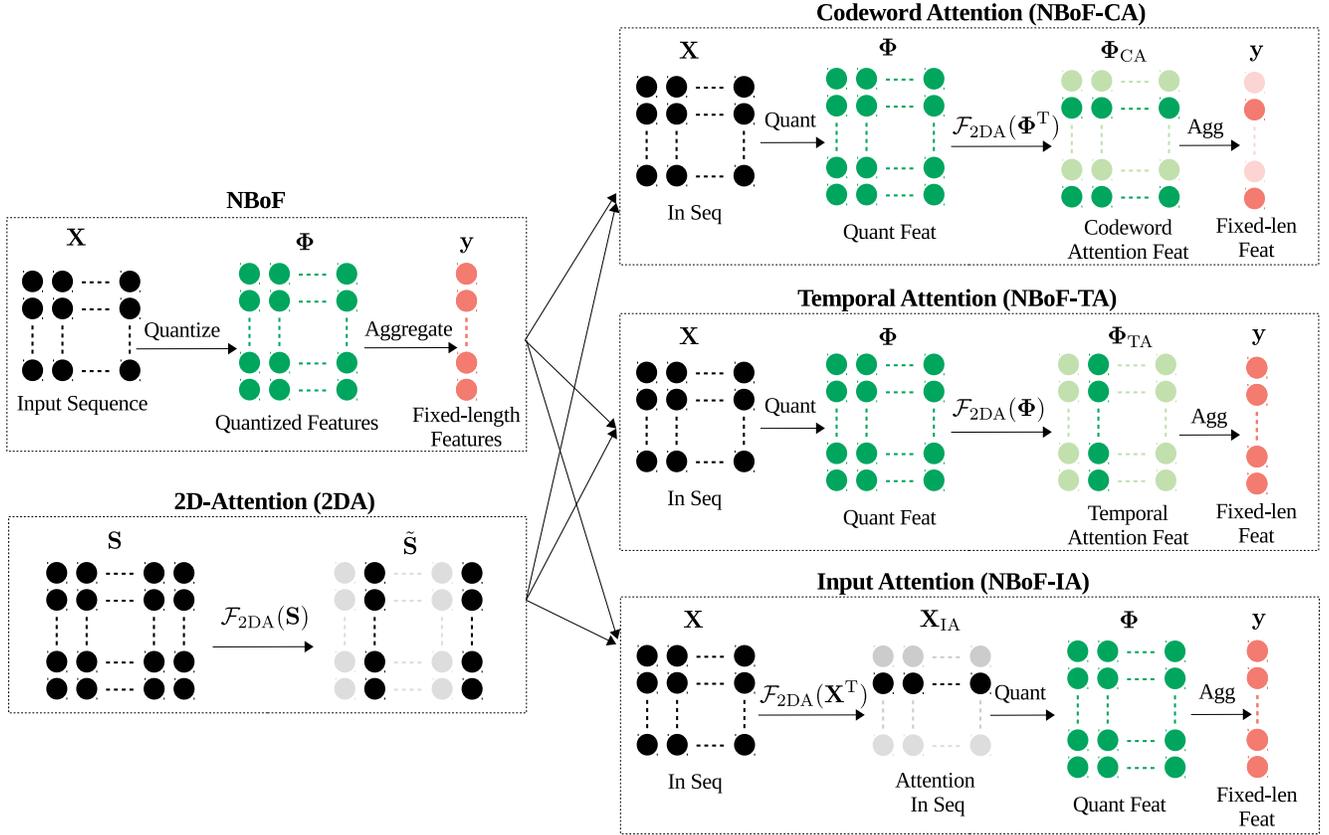
Fig. 1. Illustration of the proposed attention formulation (2DA) and different attention-based NBoF models

$$\mathbf{A} = \mathcal{G}(\mathbf{Z})$$
$$\mathbf{Z} = \mathbf{SW} \qquad (5)$$

where $\mathcal{G}(\mathbf{Z})$ denotes the soft-max function that is applied to every row of $\mathbf{Z}$. That is, every element of $\mathbf{A}$ is non-negative, and each row of $\mathbf{A}$ sums up to 1. Similar to other attention formulations [25], [6], [22], we use soft-max normalization to promote competitions between different columns of $\mathbf{Z}$.

As mentioned previously, the weight matrix $\mathbf{W}$ is used to measure the relative importance between columns of $\mathbf{S}$, which is encoded in $\mathbf{Z}$, and thus $\mathbf{A}$. In order to see this, let us denote by $\mathbf{s}_n \in \mathbb{R}^M$ and $\mathbf{z}_n \in \mathbb{R}^M$ the $n$-th column of $\mathbf{S}$ and $\mathbf{Z}$, respectively. Since $\mathbf{Z} = \mathbf{SW}$, the $n$-th column of $\mathbf{Z}$, i.e., $\mathbf{z}_n$, is calculated as the weighted combination of $N$ columns of $\mathbf{S}$, with the weight of the $n$-th column always equal to $1/N$ since the diagonal elements of $\mathbf{W}$ are fixed to $1/N$. In this way, element $z_{mn}$ (in $\mathbf{Z}$) encodes the relative importance of $s_{mn}$ (in $\mathbf{S}$) with respect to other $s_{mk}$, for $k \neq n$.

### B. Attention-based Neural Bag-of-Features

In this subsection, we will show how the proposed attention module 2DA can be used to address different limitations of the NBoF model described in Section I.

**Codeword Attention**: in the NBoF model, quantization results produced by each quantization neuron (codeword) are considered equally important for every input sequence. This property limits the feature quantization step to fully take advantage of the data-driven approach. In order to overcome this limitation, the proposed 2DA block can be applied to the quantized features to highlight or discard the outputs of certain quantization neurons. By doing so, the NBoF model is explicitly encouraged to learn a subset of specialized codewords for a given input pattern.

Particularly, given the quantized features denoted by $\mathbf{\Phi} \in \mathbb{R}^{K \times N}$ as described in Section II, we propose to apply attention to the rows of $\mathbf{\Phi}$ because the first mode of $\mathbf{\Phi}$ with dimension $K$ denotes the number of quantization neurons or codewords. Since 2DA operates on the columns of the input matrix, the attention-based quantized features is calculated as follows:

$$\mathbf{\Phi}_{\mathrm{CA}} = \mathcal{F}_{\mathrm{2DA}}(\mathbf{\Phi}^{\mathrm{T}}) \qquad (6)$$

where $\mathbf{\Phi}^{\mathrm{T}}$ denotes the transpose of $\mathbf{\Phi}$.

**Temporal Attention**: another limitation of the NBoF model lies in the aggregation step. In order to produce a fixed-length representation of the input sequence, the aggregation step in the NBoF model simply computes the mean of quantized features along the temporal mode. In this way, the NBoF model only allows equal contributions of all quantized features, disregarding the temporal information. In fact, the idea

of giving different weights to different time instances has been adopted in previous works under different formulations [10], [25]. Using our proposed 2DA formulation, it is straightforward to enable the NBoF model to attend to salient temporal information as follows:

$$\mathbf{\Phi}_{\text{TA}} = \mathcal{F}_{2\text{DA}}(\mathbf{\Phi}) \tag{7}$$

Since each column of $\mathbf{\Phi}$ contains quantized features of each time step, to obtain temporal attention-based features $\mathbf{\Phi}_{\text{TA}}$ we simply apply $\mathcal{F}_{2\text{DA}}$ to $\mathbf{\Phi}$ as in Eq. (7). $\mathbf{\Phi}_{\text{TA}}$ is then averaged along the second dimension to produce the fixed-length representation of the input sequence. Although we still perform averaging in the aggregation step, the fixed-length representation is no longer the average of the quantized features, but a weighted average. This is because each time instance (column) in $\mathbf{\Phi}_{\text{TA}}$ has been scaled by different factors via the attention mechanism.

**Input Attention**: noisy data is an inherent problem in many real-world applications. Noises might surface during the data acquisition process, such as ambient noise in audio signals or power line interference and motion artifacts in Electrocardiogram signals. In other scenarios, noises are inherent in the problem formulation since the relevance between the input sources and the targets might be unclear. For example, in stock prediction, it is intuitive to use related stocks' data, e.g., those coming from the same market sector, as the input to construct forecasting models, although some of them might be irrelevant to the movement of the target stock.

The proposed attention mechanism can also be used to filter out potential noisy series in a multivariate series as follows:

$$\mathbf{X}_{\text{IA}} = \mathcal{F}_{2\text{DA}}(\mathbf{X}^{\text{T}}) \tag{8}$$

where $\mathbf{X} \in \mathbb{R}^{D \times N}$ denotes an input sequence of $N$ steps of the NBoF model as specified in Section II. Since we would like to apply attention over the individual series (rows of $\mathbf{X}$), $\mathcal{F}_{2\text{DA}}$ is applied to the transpose of $\mathbf{X}$.

The proposed attention variants of the NBoF model are illustrated in Figure 1.

## IV. EXPERIMENTS

In this section, we provide detailed descriptions and results of our empirical analysis, which demonstrate the advantages of attention-based NBoF models proposed in Section III. Experiments were conducted in different types of sequence data, namely financial time-series in stock movement prediction problem, Electrocardiogram (ECG) and Phonocardiogram (PCG) in heart anomaly detection problems, and audio recording in music genre recognition and acoustic scene classification problems.

The experiments were conducted with the recently proposed logistic formulation of the NBoF model [26], i.e., the hyperbolic kernel was used in the quantization layer. In addition, we also experimented with the temporal variant of the NBoF model as proposed in [27] with a long-term and a short-term codebook. This variant is denoted as TNBoF. The codebook

attention, temporal attention, and input attention when applied to the NBoF model are denoted as NBoF-CA, NBoF-TA, and NBoF-IA, respectively. The corresponding attention variants for the TNBoF model are denoted as TNBoF-CA, TNBoF-TA, and TNBoF-IA. In addition to the NBoF and TNBoF models serving as the baseline models, we also evaluated RNN models using Gated Recurrent Units (GRU) [5].

### A. Financial Forecasting Experiments

Although extensively studied over the last decades, financial forecasting still remains as the most challenging tasks among time-series predictions [29]. This is due to the complex dynamics of the financial markets, which make the observed data highly stationary and noisy. For this reason, we selected the stock movement prediction problem in FI2010 dataset [2] as a representative problem in time-series forecasting. FI2010 is the largest publicly available Limit Order Book (LOB) dataset, which contains approximately $4.5$ million order events. The limit orders came from $5$ Finnish stocks traded in Helsinki Exchange (operated by NASDAQ Nordic) over $10$ business days. At each time instance, the dataset provides information (the prices and volumes) of the top $10$ levels, leading to a $40$-dimensional vector representation.

The FI2010 dataset is used to investigate the problem of mid-price movement prediction in the next $H = \{10, 20, 50\}$ order events. The mid-price at a given time instance is the average between the best buy and best sell prices. This quantity is a virtual price since no trade can happen at this particular price at the given time instance. The movement of mid-price (stationary, increasing, decreasing) reflects the dynamic of the LOB and the market, thus plays an important role in financial analysis. The dataset provides the movement labels, given the future horizon $H = \{10, 20, 50\}$. Details regarding the FI2010 dataset and LOB can be found in [2].

We followed the same experimental setup proposed in [10], which used the first 7 days for training the models and the last 3 days to test the performances. Due to the imbalanced nature of the problem, we reported averaged F1 score per movement as the main performance metric, similar to prior experiments [1], [10]. Detailed information about the training hyper-parameters and the network architectures is provided in the Appendix.

The experiment results for FI2010 are shown in Table I. In the second column of Table I, we list the performances of all models without using any convolution layers as the preprocessing layers. That is, the results in the second column of Table I are produced by architectures consisting of only the layer of interests (such as GRU, NBoF, and so on), plus the fully connected layers for generating predictions. In this setting, the GRU models outperform all variants of the NBoF model. This is expected since the NBoF model, by construction, is not designed to capture local features and long-term dependency in the input sequence. We can easily observe that this limitation can be partially overcome with the TNBoF variant, which uses two separate codebooks to capture the short-term and long-term dependency. By applying

| Models | without conv | with conv |
|---|---|---|
| *Prediction Horizon $H = 10$* | | |
| GRU [5] | $\mathbf{60.92}_{\pm00.09}$ | $62.21_{\pm00.30}$ |
| NBoF [26] | $33.07_{\pm00.66}$ | $66.34_{\pm00.60}$ |
| NBoF-CA (our) | $40.81_{\pm00.05}$ | $67.56_{\pm00.02}$ |
| NBoF-TA (our) | $40.83_{\pm00.21}$ | $\mathbf{67.98}_{\pm00.09}$ |
| TNBoF [27] | $36.66_{\pm00.51}$ | $66.74_{\pm00.36}$ |
| TNBoF-CA (our) | $45.61_{\pm00.16}$ | $67.76_{\pm00.05}$ |
| TNBoF-TA (our) | $45.97_{\pm00.15}$ | $67.88_{\pm00.13}$ |
| *Prediction Horizon $H = 20$* | | |
| GRU [5] | $\mathbf{51.61}_{\pm00.25}$ | $53.83_{\pm00.14}$ |
| NBoF [26] | $38.06_{\pm00.53}$ | $58.85_{\pm00.05}$ |
| NBoF-CA (our) | $40.08_{\pm00.07}$ | $59.31_{\pm00.44}$ |
| NBoF-TA (our) | $40.34_{\pm00.06}$ | $\mathbf{60.10}_{\pm00.03}$ |
| TNBoF [27] | $38.67_{\pm00.50}$ | $59.61_{\pm00.48}$ |
| TNBoF-CA (our) | $43.06_{\pm00.03}$ | $59.73_{\pm00.19}$ |
| TNBoF-TA (our) | $43.50_{\pm00.15}$ | $60.04_{\pm00.24}$ |
| *Prediction Horizon $H = 50$* | | |
| GRU [5] | $\mathbf{63.13}_{\pm00.19}$ | $65.93_{\pm00.03}$ |
| NBoF [26] | $48.25_{\pm00.25}$ | $68.84_{\pm02.29}$ |
| NBoF-CA (our) | $49.34_{\pm00.17}$ | $73.25_{\pm00.27}$ |
| NBoF-TA (our) | $49.21_{\pm00.16}$ | $73.02_{\pm00.04}$ |
| TNBoF [27] | $54.06_{\pm00.14}$ | $69.27_{\pm01.09}$ |
| TNBoF-CA (our) | $57.15_{\pm00.21}$ | $\mathbf{73.77}_{\pm00.37}$ |
| TNBoF-TA (our) | $57.41_{\pm00.06}$ | $73.40_{\pm00.08}$ |

| Models | adaptive scale | no adaptive scale |
|---|---|---|
| *Prediction Horizon $H = 10$* | | |
| NBoF-CA | $66.92_{\pm00.08}$ | $\mathbf{67.56}_{\pm00.02}$ |
| NBoF-TA | $67.34_{\pm00.14}$ | $\mathbf{67.98}_{\pm00.09}$ |
| TNBoF-CA | $\mathbf{67.84}_{\pm00.16}$ | $67.76_{\pm00.05}$ |
| TNBoF-TA | $67.16_{\pm00.32}$ | $\mathbf{67.88}_{\pm00.13}$ |
| *Prediction Horizon $H = 20$* | | |
| NBoF-CA | $59.25_{\pm00.14}$ | $\mathbf{59.31}_{\pm00.44}$ |
| NBoF-TA | $59.26_{\pm00.18}$ | $\mathbf{60.10}_{\pm00.03}$ |
| TNBoF-CA | $\mathbf{59.75}_{\pm00.35}$ | $59.73_{\pm00.19}$ |
| TNBoF-TA | $59.78_{\pm00.11}$ | $\mathbf{60.04}_{\pm00.24}$ |
| *Prediction Horizon $H = 50$* | | |
| NBoF-CA | $71.93_{\pm00.14}$ | $\mathbf{73.25}_{\pm00.27}$ |
| NBoF-TA | $47.89_{\pm23.87}$ | $\mathbf{73.02}_{\pm00.04}$ |
| TNBoF-CA | $71.30_{\pm00.29}$ | $\mathbf{73.77}_{\pm00.37}$ |
| TNBoF-TA | $72.32_{\pm00.05}$ | $\mathbf{73.40}_{\pm00.08}$ |

our proposed attention mechanism, performances of both the NBoF and TNBoF models are further boosted.

The third column of Table I shows the performances of all models when using two additional convolution layers as the local feature extractor, prior to applying the layer of interest. It is clear that all of the models benefit from the additional convolution layers, especially the NBoF model and its variants. In this setting, the GRU models no longer dominate the family of NBoF models. In fact, the GRU models become the worst-performing ones in the third column of Table I. Furthermore, both codebook attention (NBoF-CA, TNBoF-CA) and temporal attention (NBoF-TA, TNBoF-TA) consistently enhance the baselines' performances, making attention-based models the best-performing ones.

Here we should note that although the baseline models (NBoF, TNBoF) use the adaptive scaling step proposed in [26] to improve training stability, we did not employ this step in attention-based models. The reason stems from the fact that adaptive scaling introduces additional degrees of freedom

to the quantization step, which counteracts the constraining effects of the attention mechanism. Table II shows the performances of attention-based models on the FI2010 dataset, with and without the adaptive scaling step proposed in [26]. In most cases, the adaptive scaling step slightly degrades the performances of the attention-based models. As we will see in the next subsection, this effect is more noticeable in audio datasets.

### B. Audio Analysis Experiments

One of the important types of sequence data is audio recordings. In this subsection, we present our empirical analysis using two audio datasets, representing two different applications in audio signal analysis: music genre recognition and acoustic scene classification.

In the first application, the objective is to train an acoustic system that recognizes the genre of a short musical recording. For this purpose, we conducted experiments using the *small subset* of the FMA dataset [30], which contains 8000 tracks coming from the 8 most popular genres: *pop, instrumental, experimental, folk, rock, international, electronic*, and *hip-hop*. Each audio clip is 30s long, which is transformed to Mel-spectrogram representation with 128 frequency bands using a window of 10ms with an overlap of 2.5ms. The preprocessing step results in the input sequence having dimensions of $128 \times 640$.

In the second application, the objective is to train an acoustic system that can classify the type of environment based on its surrounding sounds. For this application, we used the TUT-UAS2018 dataset [31], which contains 8640

| Models | without conv | with conv |
|---|---|---|
| **FMA Dataset** | | |
| GRU [5] | $33.87_{\pm00.27}$ | $42.06_{\pm00.92}$ |
| NBoF [26] | $35.65_{\pm01.41}$ | $38.83_{\pm01.83}$ |
| NBoF-CA (our) | $\mathbf{38.29}_{\pm00.95}$ | $41.08_{\pm01.85}$ |
| NBoF-TA (our) | $36.79_{\pm00.29}$ | $41.46_{\pm01.64}$ |
| TNBoF [27] | $35.25_{\pm03.50}$ | $39.13_{\pm00.53}$ |
| TNBoF-CA (our) | $37.29_{\pm00.66}$ | $\mathbf{42.67}_{\pm01.23}$ |
| TNBoF-TA (our) | $36.50_{\pm01.49}$ | $42.58_{\pm00.91}$ |
| **TUT-UAS2018 Dataset** | | |
| GRU [5] | $56.83_{\pm00.78}$ | $56.89_{\pm00.93}$ |
| NBoF [26] | $52.02_{\pm00.18}$ | $55.92_{\pm01.40}$ |
| NBoF-CA (our) | $\mathbf{56.89}_{\pm00.17}$ | $\mathbf{57.68}_{\pm00.65}$ |
| NBoF-TA (our) | $56.09_{\pm00.25}$ | $57.63_{\pm00.30}$ |
| TNBoF [27] | $52.62_{\pm00.78}$ | $55.30_{\pm00.13}$ |
| TNBoF-CA (our) | $56.19_{\pm00.23}$ | $56.73_{\pm00.51}$ |
| TNBoF-TA (our) | $56.34_{\pm00.62}$ | $57.33_{\pm00.20}$ |

| Models | adaptive scale | no adaptive scale |
|---|---|---|
| **FMA Dataset** | | |
| NBoF-CA | $38.37_{\pm02.04}$ | $\mathbf{41.08}_{\pm01.85}$ |
| NBoF-TA | $34.29_{\pm05.42}$ | $\mathbf{41.46}_{\pm01.64}$ |
| TNBoF-CA | $39.96_{\pm00.66}$ | $\mathbf{42.67}_{\pm01.23}$ |
| TNBoF-TA | $40.21_{\pm00.16}$ | $\mathbf{42.58}_{\pm00.91}$ |
| **TUT-UAS2018 Dataset** | | |
| NBoF-CA | $40.57_{\pm14.36}$ | $\mathbf{57.68}_{\pm00.65}$ |
| NBoF-TA | $56.62_{\pm00.39}$ | $\mathbf{57.63}_{\pm00.30}$ |
| TNBoF-CA | $56.59_{\pm00.51}$ | $\mathbf{56.73}_{\pm00.51}$ |
| TNBoF-TA | $55.05_{\pm01.25}$ | $\mathbf{57.33}_{\pm00.20}$ |

| Models | test accuracy |
|---|---|
| **Noisy FMA Dataset** | |
| GRU [5] | $31.04_{\pm01.43}$ |
| NBoF [26] | $31.54_{\pm00.21}$ |
| NBoF-IA (our) | $36.21_{\pm01.93}$ |
| TNBoF [27] | $30.71_{\pm00.87}$ |
| TNBoF-IA (our) | $\mathbf{36.67}_{\pm00.95}$ |
| **Noisy TUT-UAS2018 Dataset** | |
| GRU [5] | $56.17_{\pm01.31}$ |
| NBoF [26] | $41.73_{\pm12.66}$ |
| NBoF-IA (our) | $\mathbf{56.79}_{\pm00.86}$ |
| TNBoF [27] | $51.48_{\pm00.85}$ |
| TNBoF-IA (our) | $56.04_{\pm00.42}$ |

audio clips recorded from 10 urban acoustic scenes: *airport, shopping_mall, metro_station, street_pedestrian, public_square, street_traffic, tram, bus, metro, park*. Similar to the FMA dataset, we also transformed each audio clip to Mel-spectrogram with 128 frequency bands using a window of 40ms with an overlap of 20ms, which results in the input sequence of size $128 \times 500$.

For both applications, we report the test accuracy as the performance metric. Experiment results on FMA and TUT-UAS2018 dataset are shown in Table III. Performances of all models with and without using convolution layers for feature extraction are presented in the second and third columns, respectively.

In the FMA dataset, we can easily observe significant improvements in all models when using additional convolution layers. Without any convolution layer, the NBoF and TNBoF models outperform the GRU model on average, however, with larger variances. The order reverses when additional convolution layers were used: the GRU model enjoys a huge benefit from the preprocessing layers, outperforming the NBoF and TNBoF models. In both scenarios, i.e., with or without convolution layers, the proposed attention block greatly enhances the baseline NBoF and TNBoF models, making them the best performing models in this task.

In the TUT-UAS2018 dataset, while adding convolution layers leads to noticeable improvements for the baseline NBoF and TNBoF, we observe no similar improvement for the GRU

model. Similar to the FMA dataset, we observe consistent performance boost in the TUT-UAS2018 dataset by incorporating the proposed attention block to the NBoF and TNBoF models.

Similar to Section IV-A, we also conducted experiments in FMA and TUT-UAS2018 datasets to analyze the effects of the adaptive scaling step proposed in [26]. The results are shown in Table IV. The results obtained from both audio analysis tasks in Table IV are consistent with what we observe from the stock movement prediction task in Table II: although the adaptive scaling step can enhance NBoF and TNBoF models as demonstrated in [26], the additional degrees of freedom introduced by this step negates the competition effects enforced by the attention mechanism, leading to performance degradation when combining both methods.

In order to evaluate how well the proposed input attention

TABLE VI
PERFORMANCE (AVERAGED F1) ON AF DATASET

| Models | F1 |
|---|---|
| GRU [5] | $76.42_{\pm 00.86}$ |
| NBoF [26] | $78.15_{\pm 00.82}$ |
| NBoF-CA (our) | $78.73_{\pm 00.71}$ |
| NBoF-TA (our) | $78.55_{\pm 00.90}$ |
| TNBoF [27] | $78.27_{\pm 01.02}$ |
| TNBoF-CA (our) | $78.71_{\pm 00.92}$ |
| TNBoF-TA (our) | $\mathbf{79.52}_{\pm 00.81}$ |

TABLE VII
PERFORMANCE (MEAN OF SENSITIVITY AND SPECIFICITY) ON PCG
DATASET. THE HIGHER, THE BETTER.

| Models | Anomaly Detection | Quality Detection |
|---|---|---|
| GRU [5] | $\mathbf{90.08}_{\pm 00.68}$ | $\mathbf{72.74}_{\pm 01.40}$ |
| NBoF [26] | $50.31_{\pm 00.42}$ | $49.69_{\pm 00.34}$ |
| NBoF-CA (our) | $88.09_{\pm 00.25}$ | $71.98_{\pm 03.00}$ |
| NBoF-TA (our) | $89.32_{\pm 01.02}$ | $72.57_{\pm 02.20}$ |
| TNBoF [27] | $54.12_{\pm 05.18}$ | $53.40_{\pm 04.21}$ |
| TNBoF-CA (our) | $88.68_{\pm 00.95}$ | $72.34_{\pm 01.32}$ |
| TNBoF-TA (our) | $88.81_{\pm 01.07}$ | $69.45_{\pm 00.89}$ |

mechanism (NBoF-IA, TNBoF-IA) tackles noisy data, we simulated contaminated audio data by adding 10 synthetic frequency bands, which are generated by adding white noise to the averaged Mel coefficients. Here we should note that in this set of experiments, we did not use any convolution layers in order to gauge how well the layers of interests are resilient to noise. The results are shown in Table V. As can be seen from Table V, when moving from the noiseless to the noisy version of FMA and TUT-UAS2018 datasets, the accuracy of NBoF and TNBoF models dropped significantly. GRU models also exhibited similar behaviors, although the performance drops are less significant as compared to the NBoF and TNBoF. By incorporating the proposed input attention block to the NBoF and TNBoF models, we were able to achieve very similar performances compared to the noiseless scenario.

*C. Medical Diagnosis Experiments*

Medical diagnosis, which plays a crucial role in ensuring human prosperity, is inherently an intricate process. The quality of the diagnosis is highly dependent on the expertise of the examiner. Since it takes several years and a great amount of resources to train human experts, medical diagnosis tools have been actively developed over the past decades to assist human examiners. In our empirical analysis using medical data, we investigated the effectiveness of the proposed models in diagnosing cardiovascular diseases using publicly available Electrocardiogram (ECG) and Phonocardiogram (PCG) signals.

The AF dataset focuses on the problem of atrial fibrillation detection from ECG recordings, which are provided as the development data (training set) in the Physionet/Computing in Cardiology Challenge 2017 [32]. The dataset contains 8528 single-lead ECG recordings lasting from 9 to 60 seconds. The objective of the challenge was to classify a given recording into one of the 4 classes: normal sinus rhythm, atrial fibrillation, alternative rhythm, and noise. We followed an experimental setup similar to [33], which evaluates a given model using 5-fold cross-validation. Additionally, the recordings were clipped or padded so that they have a constant length of 30 seconds. Since a single lead ECG recording is only a univariate sequence, it is necessary to use convolution layers as preprocessing layers to extract higher-level features,

before the NBoF or GRU layers. To tackle the imbalanced nature of the training set, we scaled the loss term associated with each class, with the factor inversely proportional to the number of samples in that class.

PCG signal is often used in ambulatory diagnosis in order to evaluate the heart hemodynamic status and detect potential cardiovascular problems. The data used in our experiments come from the training set provided in the Physionet/Computing in Cardiology Challenge 2016 [34]. The objective of the challenge is to develop an automatic classification method for the anomaly (normal versus abnormal) and quality (good versus bad) detection given a PCG recording.

Since the length of the recordings varies greatly, from 5 to 120 seconds, we generated 5s segments from the recordings for training the models; during the test phase, the models were used to classify 5s sub-segments (with 4s overlap) of a given recording, and the overall label is inferred from the averaged classification of the sub-segments. PCG signal captures the acoustic nature of the heart sound; thus, we extracted Mel-spectrogram with 24 frequency bands, using a window of 25ms with an overlap of 10ms to represent each segment. With a smaller size compared to the AF dataset, we only employed a 3-fold cross-validation protocol for this problem. Further details regarding our experimental setup in AF and PCG datasets are provided in the Appendix.

Table VI shows the averaged F1 score, a metric adopted by the database [32], of all models on the AF dataset. In Table VII, we show the anomaly and quality detection performance. The performance metric used by the database [34] is calculated as the mean of sensitivity and specificity scores. In the AF dataset, the averaged F1 scores obtained from the baseline NBoF and TNBoF models are significantly higher than the one obtained from the recurrent model. Although the improvement margins for the NBoF model are minor in Table VI, both the NBoF and TNBoF models enjoy increases in performance when using the proposed attention blocks. The consistent performance gain produced by the attention blocks can also be observed in the PCG dataset in Table VII. In this dataset, while the NBoF and TNBoF models score far below the GRU model, the attention-based models perform nearly as well as

the recurrent model.

## V. CONCLUSIONS

In this paper, we proposed 2D-Attention, a generic attention mechanism for data represented in the form of matrices. The proposed attention computation can be used in a plug-and-play manner, and can be updated jointly with other components in a computation graph. Using the proposed attention block, we further proposed three variants of the Neural Bag-of-Features model when learning with sequence data. Our extensive experiments in financial forecasting, audio analysis and medical diagnosis demonstrated that the proposed attention consistently led to performance gains for the Neural Bag-of-Features models. Since 2D-Attention is a generic attention computation method for matrices, investigating its efficacy in other neural network models is an interesting research direction in the future works.

## VI. ACKNOWLEDGEMENT

## APPENDIX

In all of our experiments, we used ADAM optimizer for stochastic optimization. Weight decay (0.0001) or max-norm constraint (4.0) was used to for regularization. In addition, dropout (0.2) was applied to the output of the layer before the classification layer. In all models, before the output layer, there is a fully-connected layer with 512 neurons. For NBoF, TNBoF and the attention models, we used 256 codewords in the quantization layer. Correspondingly, the number of units in GRU model was set to 256. Details that are specific to each experiment are provided below:

- *Financial Forecasting Experiments*: All models were trained for 80 epochs, with the initial learning rate set to 0.001. The learning rate was decreased by a factor of 0.1 at epoch 11 and 51. We followed [10] and scaled the loss term associated with each class with a factor that is inversely proportional to the number of samples of each class to counter the effect of class imbalanced. In experiments that used convolution layers as preprocessing layers, we used two 1D convolution layers, each of which has 64 filters with the filter size set to 5 and the stride set to 1. Batch normalization was used after each convolution layer, followed by the ReLU activation.
- *Audio Analysis Experiments*: The setup is similar to the financial forecasting experiments, except for the configuration of convolution layers: four 1D convolution layers with the filter size of 5 were used; the first two convolution layers have 32 filters, which are followed by a max-pooling layer to reduce the temporal dimension by half. The last two convolution layers have 64 filters. After

each convolution layer, we applied batch normalization, followed by ReLU activation.
- *Medical Diagnosis Experiments*: in both AF and PCG datasets, all models were trained for 90 epochs, with the initial learning rate set to 0.001, which was decreased to 0.0001 at epoch 11, then to 0.00001 at epoch 71. For the AF dataset, we adopted the convolution architecture proposed in [33] as the first computation block in all models. For PCG dataset, we used five 1D convolution layers with the filter size set to 3 as the preprocessing layers: the first two layers have 32 filters with strides of 1; the third layer has 64 filters with strides of 2; the fourth layer has 64 filters with strides of 1; the last layer has 128 filters with strides of 2. After each convolution layer, we applied batch normalization, followed by ReLU activation.

## REFERENCES

[1] D. T. Tran, M. Magris, J. Kanniainen, M. Gabbouj, and A. Iosifidis, "Tensor representation in high-frequency financial data for price change prediction," in *2017 IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 1–7, IEEE, 2017.

[2] A. Ntakaris, M. Magris, J. Kanniainen, M. Gabbouj, and A. Iosifidis, "Benchmark dataset for mid-price forecasting of limit order book data with machine learning methods," *Journal of Forecasting*, vol. 37, no. 8, pp. 852–866, 2018.

[3] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *2013 IEEE international conference on acoustics, speech and signal processing*, pp. 6645–6649, IEEE, 2013.

[4] D. Stowell, D. Giannoulis, E. Benetos, M. Lagrange, and M. D. Plumbley, "Detection and classification of acoustic scenes and events," *IEEE Transactions on Multimedia*, vol. 17, no. 10, pp. 1733–1746, 2015.

[5] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.

[6] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.

[7] E. Slutzky, "The summation of random causes as the source of cyclic processes," *Econometrica: Journal of the Econometric Society*, pp. 105–146, 1937.

[8] G. C. Tiao and G. E. Box, "Modeling multiple time series with applications," *journal of the American Statistical Association*, vol. 76, no. 376, pp. 802–816, 1981.

[9] D. T. Tran, M. Gabbouj, and A. Iosifidis, "Multilinear class-specific discriminant analysis," *Pattern Recognition Letters*, vol. 100, pp. 131–136, 2017.

[10] D. T. Tran, A. Iosifidis, J. Kanniainen, and M. Gabbouj, "Temporal attention-augmented bilinear network for financial time-series data analysis," *IEEE transactions on neural networks and learning systems*, vol. 30, no. 5, pp. 1407–1418, 2018.

[11] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[12] N. Passalis and A. Tefas, "Neural bag-of-features learning," *Pattern Recognition*, vol. 64, pp. 277–294, 2017.

[13] B. Zhao, H. Lu, S. Chen, J. Liu, and D. Wu, "Convolutional neural networks for time series classification," *Journal of Systems Engineering and Electronics*, vol. 28, no. 1, pp. 162–169, 2017.

[14] N. Passalis, A. Tsantekidis, A. Tefas, J. Kanniainen, M. Gabbouj, and A. Iosifidis, "Time-series classification using neural bag-of-features," in *2017 25th European Signal Processing Conference (EUSIPCO)*, pp. 301–305, IEEE, 2017.

[15] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2, pp. 2169–2178, IEEE, 2006.

[16] Y.-G. Jiang, C.-W. Ngo, and J. Yang, "Towards optimal bag-of-features for object categorization and semantic video retrieval," in *Proceedings of the 6th ACM international conference on Image and video retrieval*, pp. 494–501, 2007.

[17] M. Riley, E. Heinen, and J. Ghosh, "A text retrieval approach to content-based audio retrieval," in *Int. Symp. on Music Information Retrieval (ISMIR)*, pp. 295–300, 2008.

[18] A. Iosifidis, A. Tefas, and I. Pitas, "Multidimensional sequence classification based on fuzzy distances and discriminant analysis," *IEEE Transactions on Knowledge and Data Engineering*, vol. 25, no. 11, pp. 2564–2575, 2012.

[19] A. Iosifidis, A. Tefas, and I. Pitas, "Discriminant bag of words based representation for human action recognition," *Pattern Recognition Letters*, vol. 49, pp. 185–192, 2014.

[20] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proceedings of the seventh IEEE international conference on computer vision*, vol. 2, pp. 1150–1157, Ieee, 1999.

[21] N. Passalis and A. Tefas, "Entropy optimized feature-based bag-of-words representation for information retrieval," *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 7, pp. 1664–1677, 2016.

[22] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International conference on machine learning*, pp. 2048–2057, 2015.

[23] V. Mnih, N. Heess, A. Graves, *et al.*, "Recurrent models of visual attention," in *Advances in neural information processing systems*, pp. 2204–2212, 2014.

[24] F. Laakom, N. Passalis, J. Raitoharju, J. Nikkanen, A. Tefas, A. Iosifidis, and M. Gabbouj, "Bag of color features for color constancy," *arXiv preprint arXiv:1906.04445*, 2019.

[25] Y. Qin, D. Song, H. Chen, W. Cheng, G. Jiang, and G. Cottrell, "A dual-stage attention-based recurrent neural network for time series prediction," *arXiv preprint arXiv:1704.02971*, 2017.

[26] N. Passalis, A. Tefas, J. Kanniainen, M. Gabbouj, and A. Iosifidis, "Temporal logistic neural bag-of-features for financial time series forecasting leveraging limit order book data," *arXiv preprint arXiv:1901.08280*, 2019.

[27] N. Passalis, A. Tefas, J. Kanniainen, M. Gabbouj, and A. 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.

[28] Y. G. Cinar, H. Mirisaee, P. Goswami, E. Gaussier, A. Aït-Bachir, and V. Strijov, "Position-based content attention for time series forecasting with sequence-to-sequence rnns," in *International Conference on Neural Information Processing*, pp. 533–544, Springer, 2017.

[29] O. B. Sezer, M. U. Gudelek, and A. M. Ozbayoglu, "Financial time series forecasting with deep learning: A systematic literature review: 2005–2019," *Applied Soft Computing*, vol. 90, p. 106181, 2020.

[30] M. Defferrard, K. Benzi, P. Vandergheynst, and X. Bresson, "Fma: A dataset for music analysis," *arXiv preprint arXiv:1612.01840*, 2016.

[31] A. Mesaros, T. Heittola, and T. Virtanen, "Tut acoustic scenes 2017, evaluation dataset," Nov. 2017.

[32] G. D. Clifford, C. Liu, B. Moody, H. L. Li-wei, I. Silva, Q. Li, A. Johnson, and R. G. Mark, "Af classification from a short single lead ecg recording: the physionet/computing in cardiology challenge 2017," in *2017 Computing in Cardiology (CinC)*, pp. 1–4, IEEE, 2017.

[33] F. Andreotti, O. Carr, M. A. Pimentel, A. Mahdi, and M. De Vos, "Comparing feature-based classifiers and convolutional neural networks to detect arrhythmia from short segments of ecg," in *2017 Computing in Cardiology (CinC)*, pp. 1–4, IEEE, 2017.

[34] G. D. Clifford, C. Liu, B. Moody, D. Springer, I. Silva, Q. Li, and R. G. Mark, "Classification of normal/abnormal heart sound recordings: The physionet/computing in cardiology challenge 2016," in *2016 Computing in Cardiology Conference (CinC)*, pp. 609–612, IEEE, 2016.