

Neuro-symbolic Models for Interpretable Time Series Classification using Temporal Logic Description

Ruixuan Yan*, Tengfei Ma[†], Achille Fokoue[†], Maria Chang[†], Agung Julius*

* Department of Electrical, Computer, and Systems Engineering, Rensselaer Polytechnic Institute, Troy, NY, USA
{yanr5, juliua2}@rpi.edu

[†] IBM T.J. Watson Research Center, IBM Research, Yorktown Heights, NY, USA
{tengfei.ma1, maria.chang}@ibm.com, achille@us.ibm.com

Abstract—Most existing Time series classification (TSC) models lack interpretability and are difficult to inspect. Interpretable machine learning models can aid in discovering patterns in data as well as give easy-to-understand insights to domain specialists. In this study, we present Neuro-Symbolic Time Series Classification (NSTSC), a neuro-symbolic model that leverages signal temporal logic (STL) and neural network (NN) to accomplish TSC tasks using multi-view data representation and expresses the model as a human-readable, interpretable formula. In NSTSC, each neuron is linked to a symbolic expression, i.e., an STL (sub)formula. The output of NSTSC is thus interpretable as an STL formula akin to natural language, describing temporal and logical relations hidden in the data. We propose an NSTSC-based classifier that adopts a decision-tree approach to learn formula structures and accomplish a multiclass TSC task. The proposed smooth activation functions enable the model to be learned in an end-to-end fashion. We test NSTSC on a real-world wound healing dataset from mice and benchmark datasets from the UCR time-series repository, demonstrating that NSTSC achieves comparable performance with the state-of-the-art models. Furthermore, NSTSC can generate interpretable formulas that match domain knowledge.

I. INTRODUCTION

Time series classification (TSC) is one challenging task in machine learning (ML) and data mining [1], [2]. Time-series data (TSD) exists in many aspects of our life, such as finance [3], clinical medicine [4], etc. Numerous TSC algorithms have been developed to extract hidden patterns from complicated and long sequences of TSD [5]. The interpretability of patterns extracted from TSD is particularly important to domain experts, especially for specialists in clinical medicine, geological science, and genetic science. Enriching the learning model to be intuitive for the users to understand the rules from the model and make better decisions is a critical task. For example, the detection of seizures using data from Electroencephalograms (EEG) is a long-existing challenge, which has been tackled with various ML methods [6], [7], whereas there is no model that can simultaneously capture data patterns and describe the decisions as human-readable formulas.

Signal temporal logic (STL) is a formal language used to express or specify temporal specifications of cyber-physical systems [8], [9]. For instance, STL formulas can describe

the features of the wound healing process as “If the protein MMP.2’s value is higher than \underline{x} and lower than \bar{x} during Day 1 to Day 10, then it’s in the Hemostasis stage.” Due to its expressivity, many *temporal logic inference (TLI)* algorithms have been proposed to analyze system executions by considering an STL formula as a classifier [10], [11]. Though the above models can learn STL formulas to describe the temporal properties of data, they face a computational complexity barrier, especially for long sequences of TSD. Classical *TLI* algorithms learn STL formulas by casting it as an optimization problem, incorporating the *quantitative satisfaction function (QSF)* (will be explained in Section V) into the objective function. Due to the non-smoothness of the QSF, the parameters are learned via expensive non-gradient-based methods. Moreover, the formulas from the above methods cannot reflect the relative importance of data at various time points because they model all time points uniformly.

In this study, we propose a novel differentiable Neuro-Symbolic framework for Time Series Classification (NSTSC) by using multi-view data representations. This framework associates a component of an STL formula with a neuron in a neural network (NN), making the model connected to an STL formula and also differentiable. Simultaneously, weights are assigned to the edges to reflect the impact of related inputs on a neuron’s output, and the new form of STL is called weighted STL (wSTL). Bounded *quantitative satisfaction functions* for wSTL are proposed to reflect the truth degree of a wSTL formula over TSD. The resulting wSTL formulas can therefore reflect the relative importance of data at various time points. Besides the raw data, we augment the data representation to obtain two additional data views: *spectral representation* and *derivative representation* (details to be explained in Section IV). As a result, NSTSC can provide wSTL formulas describing the characteristics of data in each view. Moreover, by utilizing a decision tree (DT) approach, the model’s wSTL formula structures can be automatically learned and a multiclass TSC task can be accomplished. By traversing the paths of the tree classifier, individual wSTL formulas for each data class can be learned.

The contributions of this study are described as follows.

1) We propose novel weighted semantics for STL (called wSTL) using bounded *quantitative satisfaction functions* to quantify the truth degree of a wSTL formula over TSD. Any wSTL formula can be mapped to an NN, where a neuron is associated with a symbolic representation of a wSTL (sub)formula. In this way we make the model differentiable and easily optimized.

2) We develop a neuro-symbolic framework called NSTSC to classify TSD, using features that are expressed in wSTL. A decision tree approach is adapted to learn formula structures automatically and the resulting tree classifier can classify multiclass TSD, allowing the tree to provide wSTL formulas describing the characteristic patterns of each data class.

3) NSTSC is evaluated on a mouse wound-healing dataset and a set of benchmark time-series datasets from the UCR time-series archive for comparison with state-of-the-art models. NSTSC is demonstrated to achieve competitive performance as well as express discriminatory patterns as interpretable logical formulas that match with domain knowledge.

II. RELATED WORK

A. Time Series Classification Models

Numerous TSC approaches have been developed to categorize TSD from various perspectives. Shapelet Transform (ST) [12] extracts shapelets, i.e., TSD subsequences, as features and classifies data based on their resemblance to the class representatives. ROCKET, the Random Convolutional Kernel Transform [13], applies random conventional kernels to TSD to extract summary statistics and utilizes a linear classifier for feature selection. In addition to the feature-based methods, neural networks (NNs) for TSC have recently been developed, such as residual network (ResNet), and fully convolutional NN (FCN) [14]. Ensemble learning models have also emerged to classify TSD. The Hierarchical vote Collective of Transformation-based Ensembles (HCTE) is an ensemble model that classifies data using a weighted average of predictions from its base classifiers: ST, Bag-Of-Symbolic Fourier-approximation-Symbols (BOSS), Time Series Forest (TSF), and Random Interval Spectral Ensemble (RISE). Recently, HCTE was upgraded to HCTE 2.0 [15] by replacing the base estimators with Temporal Dictionary Ensemble [16], Diverse Representation Canonical Interval Forest [17], and Arsenal [15]. The Time Series Combination of Heterogeneous and Integrated Embedding Forest (CHIEF) is an ensemble TSC model comprising trees as base estimators, which consider various features and splitting criteria for classification [18]. This work is concerned with designing a neuro-symbolic framework that can not only classify TSD but also express the model as an interpretable formula for better decision-making.

B. Signal Temporal Logic Inference

A large body of STL inference works on TSD analysis have been developed during the last decades. The majority of them infer a formula by incorporating the *quantitative satisfaction function* into the objective function and solving an optimization problem. A parametric STL (pSTL) was proposed

in [19] to find parameters for STL formulas. A CensusSTL describing the number of agents satisfying an STL formula in different groups was proposed in [20], accompanied by a CensusSTL inference algorithm. These methods have high computation costs due to the non-smooth objective functions and large search space. In this study, we introduce importance weights into STL to obtain differentiable *quantitative satisfaction functions* and develop a neuro-symbolic model that can learn the parameters of STL formulas in an end-to-end fashion.

C. Neuro-Symbolic Models

In recent years, neuro-symbolic models have been extensively investigated in various disciplines [21]–[23]. Many-valued logics, such as fuzzy logic, have strong connections with wSTL and have been applied to time series forecasting [24]. However, the model has no relation to neural networks. Neuro-symbolic models seek to combine the benefits of robust learning in neural networks with the benefits of interpretability and reasoning in symbolic representations. Markov logic networks share a similar aim, using Markov networks to model logical formulas and to handle inference under uncertainty. Logical Neural Networks (LNN) [25] takes this a step further by using neural networks to model formulas in real-valued logics. LNNs have been applied to bi-directional inference and theorem proving tasks, as well as question answering over knowledge bases [26]. Recently, a rule-based representation learner (RRL) was developed in [21] to automatically learn data representation and classification as interpretable rules, where a gradient crafting approach is adapted to allow for continuous learning. Nevertheless, these models do not apply to TSD that has no grounding truth value.

III. PRELIMINARIES

In this study, TSD refers to p -dimensional data that evolves in discrete time and is denoted as $x = \{x(0), x(1), \dots, x(K-1)\}$, where $x(k) \in \mathbb{R}^p, \forall k < K, k$ is the time index, $K \in \mathbb{Z}_+$ is the duration of x , and $p \in \mathbb{Z}_+$. The TSC task is performed on a labeled dataset $D = \{(x_i, y_i)\}_{i=1}^N$, where x_i is the i -th data in the dataset ($x_i = \{x_i(0), \dots, x_i(K-1)\}, x_i(k) \in \mathbb{R}^p$), and y_i is the label of x_i . The total number of classes in the dataset is denoted as C , meaning $y_i \in \mathcal{C} = \{1, 2, \dots, C\}$.

A. Weighted Signal Temporal Logic (wSTL)

The concept of wSTL was first introduced in [27]. In this study, we propose novel bounded quantitative semantics for wSTL based on the notion of truth degree.

Definition 1. The syntax of wSTL is defined as [27]

$$\phi := \top | \pi | \neg \phi | \phi_1^{w_1} \wedge \phi_2^{w_2} | \phi_1^{w_1} \vee \phi_2^{w_2} | \Diamond_{[k_1, k_2]}^w \phi | \Box_{[k_1, k_2]}^w \phi, \quad (1)$$

where \top is Boolean TRUE, π is an atomic predicate defined as $\pi := f(x) \geq 0$. E.g., if $f(x) = \mathbf{a}^T x - u, \|\mathbf{a}\| = 1, u \in \mathbb{R}$, then π represents a half space and \mathbf{a}, u are parameters to learn. The symbols \neg, \wedge , and \vee denote logical negation, conjunction, and disjunction, respectively. The temporal operators \Diamond and \Box mean “Eventually” and “Always”, respectively. We assume that $k_1 \leq k_2 < K$. The (unweighted) expression $\Diamond_{[k_1, k_2]} \phi$

means the formula ϕ is satisfied somewhere within the time interval $[k_1, k_2]$, while $\Box_{[k_1, k_2]}\phi$ means the formula ϕ is satisfied everywhere within the time interval $[k_1, k_2]$. See the example below for further explanation. w_1 and w_2 are nonnegative weights associated with ϕ_1 and ϕ_2 in the \wedge and \vee operators, and $\mathbf{w} = [w_{k_1}, w_{k_1+1}, \dots, w_{k_2}]^T \in \mathbb{R}_{\geq 0}^{k_2-k_1+1}$ denotes nonnegative weights associated with \Box and \Diamond operators, with $w_{k'}$ being the weight assigned to time $k' \in [k_1, k_2]$.

Example 1. To help understand the temporal logic semantics, we illustrate an example of the unweighted version of STL in Fig. 1. Two time series data, $\text{Red}(k)$ and $\text{Blue}(k)$, are shown. The green band represents a logical predicate $\pi : (x \leq 4) \wedge (x \geq 3)$. An unweighted STL formula $\phi_1 \triangleq \Diamond_{[3,4]} \pi$ means the logical predicate π is satisfied somewhere in the time interval $[3,4]$. Note that $\text{Red}(k)$ satisfies ϕ_1 while $\text{Blue}(k)$ does not. We can thus use ϕ_1 as a classifier between the two TSD. Similarly, an unweighted STL formula $\phi_2 \triangleq \Box_{[5,6]} \pi$ means the logical predicate π is satisfied everywhere in the time interval $[5,6]$. We see that $\text{Blue}(k)$ satisfies ϕ_2 while $\text{Red}(k)$ does not, implying that ϕ_2 would also be a good classifier. The weights of the formulas, as we will explain later, are used to describe the relative importance of the time indices in the temporal operations or the clauses in the logical operations.

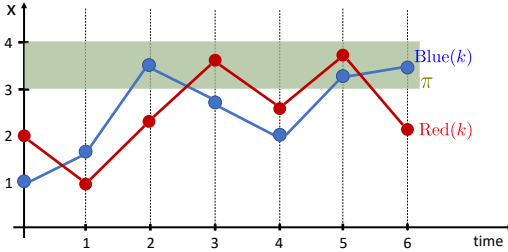


Fig. 1. Illustration of the unweighted version of the STL semantics.

IV. MULTI-VIEW DATA REPRESENTATION & FEATURE EXTRACTION

Model Structure Overview: The overall structure of the TSC model is shown in Fig. 2. The raw TSD is first transformed into multi-view data representations, which are then used to extract the interval features. The multi-view data representations and the extracted interval features are given as inputs to the NSTSC. NSTSC relies on constructing a decision tree to classify the input data, where Node i in the tree is a binary classifier designed as a neuro-symbolic model characterizing a wSTL formula ϕ_{Bi} , which describes the characteristics of data and features in each view. Traversing the path from the root node to a leaf node gives the prediction of data and a formula for the data class in that leaf node.

A. Multi-View Data Representation

1) **Raw Representation:** The original TSD is considered as the raw representation, from which spectral and derivative representations are collected.

2) **Spectral Representation:** Studies have shown classifying TSD from the frequency domain can reveal discriminatory features invisible in the raw data [28]. Spectral data representation is obtained via fast Fourier transform (FFT) as follows:

$$\begin{aligned} \hat{x}(\omega) &= \sum_{k=0}^{K-1} x(k)(e^{-i\frac{2\pi}{K}\omega k}), \quad \omega = 0, \dots, K-1, \\ &= \hat{a}(\omega) + i\hat{b}(\omega), \end{aligned} \quad (2)$$

where i is the imaginary unit. The power spectrum is then used to express the spectral representation, i.e., $\hat{x} = \{\hat{x}(0), \dots, \hat{x}(K-1)\}$, where $\hat{x}(\omega) = \sqrt{\hat{a}^2(\omega) + \hat{b}^2(\omega)}$.

3) **Derivative Representation:** For derivative representation, we adopt the same derivative transformation technique as [29], which has shown the efficacy of derivative representation in TSC. The derivative representation is obtained via computing the first-order difference, i.e., $\tilde{x}(k) = x(k+1) - x(k)$.

B. Interval Feature Extraction

Discriminatory interval features have been demonstrated to be effective in classifying TSD [28]. We split the TSD into intervals with length \mathcal{I} and extract features from each interval. Specifically, given an aggregation function $f_A(\cdot)$ and an interval $[k \cdot \mathcal{I}, (k+1) \cdot \mathcal{I} - 1]$, $f_A(\cdot)$ is applied to the interval data $\{x(k \cdot \mathcal{I}), \dots, x((k+1) \cdot \mathcal{I} - 1)\}$ to obtain an interval feature $x_{IF}(k) = f_A(x, k \cdot \mathcal{I}, (k+1) \cdot \mathcal{I} - 1)$. For example, if $f_A(\cdot) = \text{mean}(\cdot)$, $k = 0$, then the mean of the data within the interval $[0, \mathcal{I} - 1]$ is denoted as $x_{IF}(0) = \text{mean}(\{x(0), \dots, x(\mathcal{I} - 1)\})$. By applying $f_A(\cdot)$ to the entire TSD, we can extract a sequence of interval features with length M , i.e., $x_{IF} = \{x_{IF}(0), \dots, x_{IF}(M-1)\}$, where $M = \lceil \frac{K}{\mathcal{I}} \rceil$. For each data representation, the same interval feature extraction technique is applied. The interval feature of \hat{x} (resp. \tilde{x}) is denoted as \hat{x}_{IF} (resp. \tilde{x}_{IF}). The data representations and their corresponding interval features are given as inputs to the classifier, which learns a ϕ describing the characteristics of data and interval features in each view. Please see Fig. 3 for an illustration.

V. NEURO-SYMBOLIC TIME SERIES CLASSIFICATION MODULE

The neuro-symbolic time series classification module relies on constructing a tree phase classifier (TPC) from the node phase classifiers (NPCs) for multiclass classification and formula structure learning. For clarity of the presentation, the variable notations of the NSTSC module is shown in Table 1.

A. Node Phase Classifier (NPC)

An NPC aims to design a neural network for a wSTL formula ϕ , where each neuron is linked to a symbolic expression in a wSTL formula. This section first presents the activation functions (AFs), i.e., quantitative satisfaction functions, for the components in (1) such that every neuron can be mapped to a symbolic interpretation, then presents the design of an NPC using the AFs, and finally discusses the learning of an NPC.

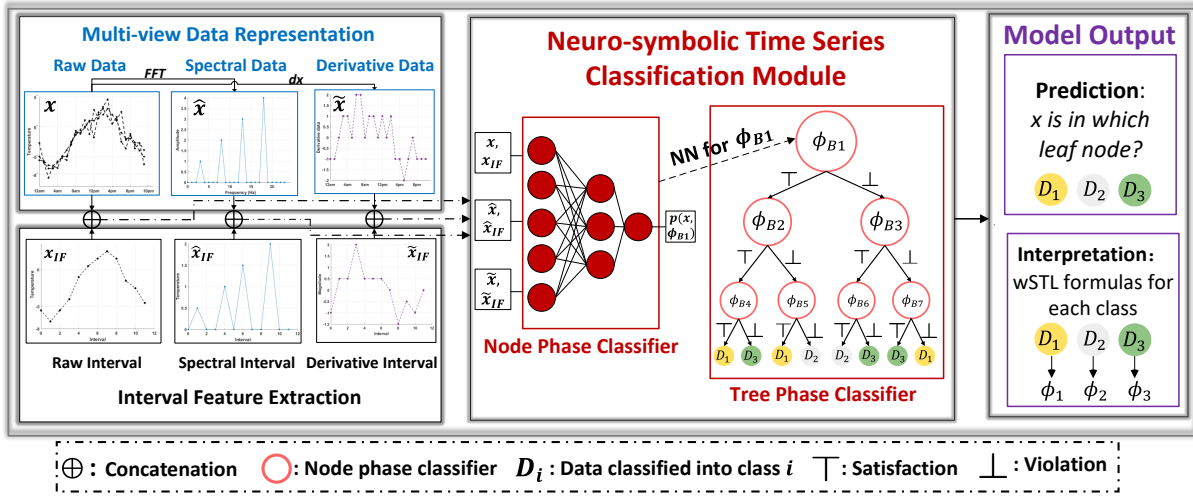


Fig. 2. The overall model structure of NSTSC, where the multi-view data and the extracted interval features are given as inputs to NSTSC, and the classification task is accomplished by learning a tree phase classifier. The leaf nodes represent data that are classified into a particular class. The tree phase classifier gives a wSTL formula describing each data class, e.g., the formula for data with label 2 is $\phi_2 = (\phi_{B1} \wedge \neg\phi_{B2} \wedge \neg\phi_{B5}) \vee (\neg\phi_{B1} \wedge \phi_{B3} \wedge \phi_{B6})$.

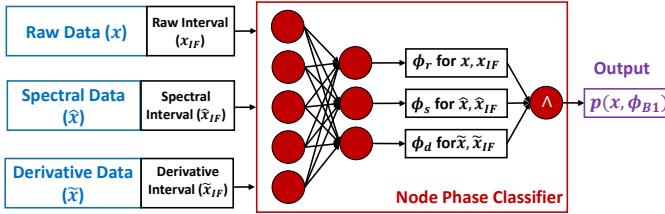


Fig. 3. Illustration of the overall workflow of the node phase classifier for ϕ_{B1} , which is in the form of $\phi_{B1} = \phi_r^{w_r} \wedge \phi_s^{w_s} \wedge \phi_d^{w_d}$. ϕ_r , ϕ_s , and ϕ_d describes the properties of data and features in the raw, spectral and derivative representations, respectively, and w_r , w_s , w_d are the corresponding weights. The node phase classifier's output is the truth degree of ϕ_{B1} over x , denoted as $p(x, \phi_{B1})$.

TABLE I
VARIABLE NOTATIONS OF THE NSTSC MODULE.

Variable notation	Representation
$x / \hat{x} / \tilde{x}$	raw / spectral / derivative data
$x_{IF} / \hat{x}_{IF} / \tilde{x}_{IF}$	raw / spectral / derivative interval features
π	atomic predicate
$\neg / \wedge / \vee$	negation / conjunction / disjunction
\diamond / \square	"Eventually" / "Always"
$p(x, \phi, k)$	truth degree of ϕ over x at k
ϕ_r, ϕ_s, ϕ_d	subformula describing the raw, spectral, derivative data
\mathcal{B}	formula structure base
D^n	data assigned to Node n
D_c^n	data in Node n with the encoded label for Class c
ϕ_B^c	formula from an NPC for class c with structure $\phi_B \in \mathcal{B}$

Definition 2. To describe the truth degree, i.e., the degree of satisfaction of ϕ over x at k , denoted as $p(x, \phi, k)$, we define the bounded quantitative satisfaction functions for wSTL as

$$\begin{aligned}
 p(x, \pi, k) &= g(f(x(k))), \\
 p(x, \neg\phi, k) &= 1 - p(x, \phi, k), \\
 p(x, \phi_1^{w_1} \wedge \phi_2^{w_2}, k) &= \otimes^{\wedge}([w_j, p(x, \phi_j, k)]_{j=1,2}), \\
 p(x, \phi_1^{w_1} \vee \phi_2^{w_2}, k) &= \oplus^{\vee}([w_j, p(x, \phi_j, k)]_{j=1,2}),
 \end{aligned} \tag{3}$$

$$\begin{aligned}
 p(x, \diamond^w \phi, k) &= \oplus^{\diamond}([w_{k'}, p(x, \phi, k + k')]_{k' \in [k_1, k_2]}), \\
 p(x, \square^w \phi, k) &= \otimes^{\square}([w_{k'}, p(x, \phi, k + k')]_{k' \in [k_1, k_2]}),
 \end{aligned}$$

where $\pi := f(x) \geq 0$, $g(\cdot)$, $\otimes^{\wedge}(\cdot)$, $\oplus^{\vee}(\cdot)$, $\otimes^{\square}(\cdot)$, $\oplus^{\diamond}(\cdot)$ are activation functions for predicates and \wedge , \vee , \square , \diamond operators, respectively. Note that (3) returns a value between 0 and 1 that represents the truth degree, different from the traditional quantitative satisfaction functions in, e.g., [10] that return values in \mathbb{R} . We denote $p(x, \phi, 0)$ as $p(x, \phi)$ for simplicity.

1) Design of Activation Functions for wSTL: The bounded quantitative satisfaction functions should be such that $p(x, \phi, k)$ close to 1 (resp. 0) indicates that x robustly satisfies (resp. violates) ϕ at k . In this study, we define $g(\cdot)$ in (3) as

$$g(r) = \text{sigmoid}(r) \triangleq \frac{1}{1 + e^{-r}}. \tag{4}$$

Note that $g(r)$ is differentiable everywhere with nonnegative derivative for any $r \in \mathbb{R}$. Essentially, temporal operators \square (resp. \diamond) can be expressed as a sequence of logical operators \wedge (resp. \vee). For instance, $\square_{[k_1, k_2]}^w \phi$ can be expressed as

$$\square_{[k_1, k_2]}^w \phi = \phi_{k_1}^{w_{k_1}} \wedge \phi_{k_1+1}^{w_{k_1+1}} \wedge \dots \wedge \phi_{k_2}^{w_{k_2}}, \tag{5}$$

where $\phi_{k'}$ denotes the formula ϕ evaluated on data instance $x(k')$ for $k' \in [k_1, k_2]$. As a result, our task is to design activation functions (AFs) for the logical operators \wedge and \vee . The AF for the \vee operator can be derived from the AF for the \wedge operator using the De Morgan's law, $\neg(\phi_1 \vee \phi_2) = (\neg\phi_1) \wedge (\neg\phi_2)$.

The AF for the \wedge operator is defined as [25]:

$$\otimes^{\wedge}([w_j, p(x, \phi_j, k)]_{j=1,2}) \triangleq h(\beta - \sum_{j=1}^2 \bar{w}_j (1 - p(x, \phi_j, k))), \tag{6}$$

where $\bar{w}_j = w_j / (w_1 + w_2)$ is the normalized weight, w_j and β are parameters to be learned, and $h(z) \triangleq \max\{0, \min\{z, 1\}\}$ is introduced to clamp the truth degree to the range $[0, 1]$.

Algorithm 1 Node Formula Learning Algorithm - *NFLA()*

Input: \mathcal{D}_c^n : TSD assigned to Node n with the encoded label for class c , ϕ_B : a basic formula in \mathcal{B} , I_t : maximum number of iterations

Output: Learned wSTL formula ϕ_B for an NPC

- 1: Construct an NSTSC based on the structure of ϕ_B , and initialize the parameters $\mathbf{w}, \beta, \mathbf{a}, u$ in ϕ_B
- 2: **for** $t = 1, 2, \dots, I_t$ **do**
- 3: Select a mini-batch data $\mathcal{D}_c^{n,t}$ from \mathcal{D}_c^n
- 4: Run forward-propagation for NSTSC to compute the truth degree of ϕ_B over $\mathcal{D}_c^{n,t}$ using (4) - (9)
- 5: Compute the loss at the current iteration using (11)
- 6: Run back-propagation to update the parameters including $\mathbf{w}, \beta, \mathbf{a}, u$ in ϕ_B
- 7: **end for**
- 8: **return** ϕ_B

the traditional NNs via back-propagation, which is described as Algorithm 1. The forward propagation computes the truth degree of ϕ_B over x , and the backward propagation updates the parameters in ϕ_B . When back propagating to the neurons for the operators, \mathbf{w} and β in (6)-(9) are updated, and when back propagating to the neurons for the predicates π , the parameters defining $f(x)$ are updated. For instance, if $f(x) = \mathbf{a}^T x - u$, then the parameters \mathbf{a} and u are updated.

B. Tree Phase Classifier (TPC)

While the expressiveness of $\phi_B \in \mathcal{B}$ is limited, combining multiple basic formulas could enrich the expressiveness. An NPC can complete a binary classification task. Multiclass TSC can be accomplished by constructing a decision tree (DT) [30] from the NPCs as a tree phase classifier (TPC). The procedures for constructing the TPC are presented in Algorithm 2, from which we could learn a formula ϕ_c for every data class c . Initially, we check if Node n is the root node. If Node n is the root node, then we set $\phi_{pt} = \emptyset$ (line 2). Next, we check if the stop condition is satisfied, e.g. $h \leq 5$ (line 4 - 6). If the stop condition is satisfied, then the tree up to the current node is returned as the tree classifier \mathcal{T} . Otherwise, we will train the node classifier for Node n . For each class $c \in \mathcal{C}$, we encode the labels of data in D^n for class c and obtain \mathcal{D}_c^n (line 8). Then for each $\phi_B \in \mathcal{B}$, we use the given ϕ_B and the encoded dataset \mathcal{D}_c^n to train an NPC for class c via Algorithm 1. By doing so, we learn a formula ϕ_B^c that is with formula structure ϕ_B and can classify data with label c and the remaining data in D^n . Next, we evaluate the performance of ϕ_B^c using the branching criterion defined as the standard Gini index, i.e.,

$$J(D^n, \phi_B^c) = \frac{|D_+^n(\phi_B^c)|}{|D^n|} G_I(D_+^n(\phi_B^c)) + \frac{|D_-^n(\phi_B^c)|}{|D^n|} G_I(D_-^n(\phi_B^c)), \quad (12)$$

where $D_+^n(\phi_B^c)$ denotes the data in D^n satisfying ϕ_B^c , and $D_-^n(\phi_B^c)$ denotes the data in D^n violating ϕ_B^c , and

$$G_I(D_+^n(\phi_B^c)) = 1 - \sum_{m=1}^C (p(D_+^n(\phi_B^c), m))^2,$$

$$p(D_+^n(\phi_B^c), m) = \frac{|\{i | (x_i, y_i) \in D_+^n(\phi_B^c), y_i = m\}|}{|D_+^n(\phi_B^c)|}$$

denotes the fraction of data in $D_+^n(\phi_B^c)$ with the original label m ,

$$G_I(D_-^n(\phi_B^c)) = 1 - \sum_{m=1}^C (p(D_-^n(\phi_B^c), m))^2,$$

$$p(D_-^n(\phi_B^c), m) = \frac{|\{i | (x_i, y_i) \in D_-^n(\phi_B^c), y_i = m\}|}{|D_-^n(\phi_B^c)|}$$

denotes the fraction of data in $D_-^n(\phi_B^c)$ with original label m (line 10). By repeating this step for each $\phi_B \in \mathcal{B}$ and each data class c , we could learn a set of formulas with different structures for each class c , from which we choose the formula ϕ_B^c with the best criterion as $\phi^{n,*}$ that represents the optimal formula at Node n (line 13). Thus a greedy search strategy is utilized to learn the formula at a node. With $\phi^{n,*}$, D^n is partitioned into $D_+^n(\phi^{n,*})$ satisfying $\phi^{n,*}$ and $D_-^n(\phi^{n,*})$ violating $\phi^{n,*}$, which is denoted as $partition(D^n, \phi_{pt} \wedge \phi^{n,*})$ in line 14. Next, we generate two child nodes of n , denoted as Node $n+1$ and $n+2$, and distribute the data $D_+^n(\phi^{n,*})$ and $D_-^n(\phi^{n,*})$ to Nodes $n+1$ and $n+2$, respectively (line 15). The DTFL algorithm is then performed on the two child nodes (line 16-17), where $n.left$ and $n.right$ denote the left and right child node of Node n , respectively. The above procedures proceed until the stop conditions are satisfied. The conversion from a TPC to a wSTL formula ϕ_c that describes the data with label c can be completed by a similar approach to [30]. The details of the conversion are presented in Appendix VIII-A.

VI. EXPERIMENTS

We evaluate NSTSC using two groups of datasets. One group is a set of real-world wound healing data from experiments on mice, and another is a set of datasets from the UCR TSD archive [31]. The experiments are run using the AdamW optimizer in Pytorch (1.10.2) on a macOS 11.4 system with a Quad-Core CPU (i7, 2.9GHz) and 16GB RAM. Our code is available at <https://tinyurl.com/NSTSC>.

A. Wound Healing Dataset

The wound healing data comes from real-world mice experiments. The task is to analyze and predict wound healing stages so that proper interventions can be done to accelerate the healing process.

1) **Dataset Description:** The wound healing process is divided into stages of hemostasis (Hem), inflammation (Inf), proliferation (Pro), and maturation (Mat). The dataset contains 67 mice data, each with two recording time points corresponding to the first day and the date of a healing stage. Every data has four features representing four protein levels: MMP.2 (x^1), IL.6 (x^2), PLGF.2 (x^3), and VEGF (x^4). We choose the following baseline models for comparison, including FCN, ResNet, and multi-layer perceptron (MLP) for TSC [14].

Algorithm 2 Decision Tree Formula Learning - $DTFL()$

Input: ϕ_{pt} : formula associated with the current path; n : the node index; h : the current depth; \mathcal{B} : formula structure base, $stop$: stop condition; J : branching criterion; D^n : TSD at Node n

Output: The tree phase classifier \mathcal{T} for multiclass TSC

```

1: if  $n = 0$  then
2:   Set  $\phi_{pt} = \emptyset$ 
3: end if
4: if  $stop(\phi_{pt}, n, h, D^n)$  then
5:   Return  $\mathcal{T}$  as the multiclass TPC
6: end if
7: for  $c \in \mathcal{C} = \{1, 2, \dots, C\}$  do
8:   Encode labels of data in  $D^n$  for class  $c$  and obtain  $D_c^n$ 
9:   for  $\phi_B \in \mathcal{B}$  do
10:    Run Algorithm 1 using  $\phi_B$  and  $D_c^n$  to learn  $\phi_B^c$ ,
    compute  $J(D^n, \phi_B^c)$  using (12)
11:   end for
12: end for
13:  $\phi^{n,*} = \arg \min_{\phi_B^c} J(D^n, \phi_B^c)$ 
14:  $n \leftarrow \phi^{n,*}, D_+^n(\phi^{n,*}), D_-^n(\phi^{n,*}) \leftarrow partition(D^n, \phi_{pt} \wedge \phi^{n,*})$ 
15:  $D^{n+1} \leftarrow D_+^n(\phi^{n,*}), D^{n+2} \leftarrow D_-^n(\phi^{n,*})$ 
16:  $n.left \leftarrow DTFL(\phi_{pt} \wedge \phi^{n,*}, n + 1, h + 1, \mathcal{B}, stop, J, D^{n+1})$ 
17:  $n.right \leftarrow DTFL(\phi_{pt} \wedge \neg \phi^{n,*}, n + 2, h + 1, \mathcal{B}, stop, J, D^{n+2})$ 

```

2) **Experiment Setup and Results:** As data is given as four proteins, we consider using a predicate in the form of $\pi_j := a^j x^j - u^j \geq 0$ for each protein, where x^j denotes the j -th protein feature, and a^j and u^j are parameters to learn. The interval features are omitted for the wound healing dataset as the sequence is short. The formula structure is chosen from $\phi = \phi_1^{w_1} \wedge \phi_2^{w_2}$ and $\phi = \phi_1^{w_1} \vee \phi_2^{w_2}$, where ϕ_1 describes data on the first day, and ϕ_2 describes data on the date of a specific stage. The task for the wound healing dataset is to predict the healing stage for a mouse given its wound data. A stratified 5-fold cross-validation approach is used for data split. The parameters a^j, u^j are initialized as 1×10^{-5} , and w are initialized as random numbers from a uniform distribution in the range of $[0, 1)$, β are initialized as 1. The maximum epoch number is 100, and the learning rate (LR) is 0.1. The mean and standard deviation (std) of accuracy, recall, precision, and F1 score are shown in Table II. We can observe that NSTSC achieves better results than the baseline models regarding higher classification accuracy and F1 score.

3) **Interpretability Analysis:** Studies have shown that MMP.2 level will increase from the Hem stage to the Pro stage [32]. The formula learned by NSTSC for the Hem stage is $\phi_h = (\phi_1^{w_1})^{0.010} \wedge (((x^1 \leq 12.61)^{0.96} \wedge (x^2 \leq 0)^{0.030} \wedge (x^3 \geq 0.9200)^{0.010})^{0.99} \vee (x^4 \geq 0)^{0.010})^{0.99}$, and the formula learned for the Pro stage is $\phi_p = (\phi_1^{w_1})^{0.080} \wedge (((x^1 \geq 17.82)^{0.57} \wedge (x^2 \geq 0)^{0.39} \wedge (x^3 \leq 0.03000)^{0.030})^{0.97} \vee (x^4 \geq 0)^{0.020})^{0.92}$,

where ϕ_h^1 and ϕ_p^1 describe data on the first day, and the value of the thresholds in the predicates are all learned by the model. If we extract the important subformulas, ϕ_h reads as “MMP.2 is below 12.61 at $k = 2$ ”, and ϕ_p reads as “MMP.2 is above 17.82 at $k = 2$ ”. This identifies “MMP.2” increases from the Hem stage to the Pro stage, matching the domain knowledge of MMP.2’s increase in [32]. Note that the value 12.61 and 17.82 are parameters u of atomic predicates π in (1) that are learned by training NSTSC. Also, we visualize “MMP.2” level at $k = 2$ and the decision rules of ϕ_h ($l \leq 12.61$) and ϕ_p ($l \geq 17.82$) in Fig. 5, where x-axis denotes the index of mice belonging to a particular stage, and y-axis denotes the “MMP.2” level (l). We can observe “MMP.2” of most mice in the Hem stage is below 12.61, and “MMP.2” of most mice in the Pro stage is above 17.82, which matches the patterns described by ϕ_h and ϕ_p . This demonstrates NSTSC can provide interpretable rules that match domain knowledge and can assist clinical experts in making better decisions on the wound healing stage prediction. Although the wound healing data has only two time points, NSTSC is proven to provide clinical experts with practically useful rules that are human-readable and easy to understand. In addition, we use the following UCR dataset to show that NSTSC also works for long sequences of TSD.

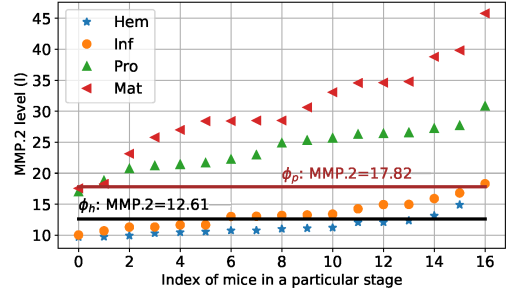


Fig. 5. Expression of “MMP.2” for four stages and decision rules in ϕ_h (black line) and ϕ_p (brown line).

B. UCR Time-Series Data

We implement NSTSC on the UCR TSD archive [31] for the second experiment. The standard split of training data and test data is utilized here. The $f(x)$ in the atomic predicates are set in the form of $f(x) = a^T x - u, \|a\| = 1$. The parameters defining the atomic predicates are initialized as 1×10^{-5} . The weights for operators w are initialized as random numbers from a uniform distribution in the range of $[0, 1)$, and β are initialized as 1, and the LR is set as 0.1. The number of training epochs is set as 100. We compare NSTSC with the following state-of-the-art (SOTA) models: FCN and ResNet [14], ST [12], BOSS [33], HCTE [28], HCTE2 [15], ROCKET [13], and CHIEF [18]. 85 commonly used datasets in the above SOTA works are selected for comparison purposes.

1) **Experimental Results:** In the experiment, the stop condition depends on the maximal depth of the DT. The maximal tree depth is five for the binary datasets, seven for the multiclass datasets with less than 10 classes, and nine

TABLE II
PERFORMANCE OF ALL TSC MODELS ON THE WOUND HEALING DATA.

Model	Accuracy		Recall		Precision		F1	
	mean	std	mean	std	mean	std	mean	std
FCN	76.99	5.03	77.66	6.63	80.68	10.16	0.7558	0.0741
ResNet	57.47	10.54	61.00	12.18	55.30	16.61	0.5449	0.1285
MLP	23.29	0.92	20.00	0	4.66	0.18	0.0755	0.0024
NSTSC	83.85	9.42	84.17	8.99	86.33	7.89	0.8327	0.1001

TABLE III
AVERAGE CLASSIFICATION ACCURACY (%) OF NSTSC AND SOTA MODELS ON THE UCR TIME-SERIES DATA ARCHIVE.

Model	FCN	ResNet	ST	BOSS	HCTE	HCTE2	ROCKET	CHIEF	NSTSC
Avg accuracy	80.91	82.47	82.23	81.12	84.71	86.12	85.07	84.78	84.66
No. rank 1st	12	10	9	11	12	25	15	20	24

TABLE IV
AVERAGE CLASSIFICATION ACCURACY (%) OF NSTSC AND SOTA MODELS ON THE BINARY DATASETS IN THE UCR ARCHIVE.

Model	FCN	ResNet	ST	BOSS	HCTE	HCTE2	ROCKET	CHIEF	NSTSC
Avg accuracy	86.11	87.57	86.43	85.98	88.09	89.49	88.83	88.67	89.07
No. rank 1st	6	4	4	7	5	6	7	8	11

TABLE V
CLASSIFICATION ACCURACY (%) OF THE ABLATION STUDY ON THE UCR TIME-SERIES DATA ARCHIVE.

Ablation Study Setting	No raw data	No spectral data	No derivative data	Only raw data
Avg accuracy	79.13	79.61	79.85	79.36

for the other datasets. Similar to [34], NSTSC adopts the mean, variance, maximum, minimum, interquartile range, and slope aggregation function to extract interval features. The number of intervals is set as 20, and the interval length \mathcal{I} is determined correspondingly. We report the accuracy on the test data using the model with the best criterion on the test data. The number of first place rank among all systems (No. rank 1st) and the average accuracy for all the datasets and for only binary datasets are shown in Table III and Table IV, respectively. Table III shows that NSTSC achieves an average accuracy of 84.66%, which is only lower than HCTE, HCTE2, CHIEF, and ROCKET. However, the gap between NSTSC and the other models is at most 1.46%. Moreover, NSTSC is ranked first more often than the other systems except for HCTE2. Table IV shows that NSTSC’s accuracy is 0.42% lower than HCTE2, but NSTSC ranks first more often than the other systems, which can be attributed to the logical properties of wSTL. Essentially, wSTL formulas classify data into two classes, meaning NSTSC is conceptually a binary classifier. Hence NSTSC achieves a better performance on the binary tasks. The DT-based method of constructing a multiclass classifier may impact the performance of the model on multiclass classification tasks. While the above results show NSTSC’s potential for TSC, its major strength is providing easily understandable and readable logical formulas to users, which is especially meaningful in practical applications. However, the other models do not possess such attributes.

2) **Interpretability Analysis:** The “BeetleFly” dataset in the UCR TSD archive is selected to demonstrate the interpretabil-

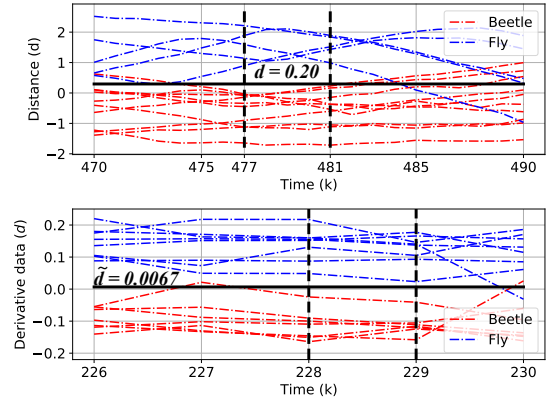


Fig. 6. Raw and derivative data of the “BeetleFly” dataset and decision rules in ϕ^{r2} (top) and ϕ^{d2} (bottom).

ity, which includes Beetle and Fly classes. The formula for the Fly data is $\phi = (\phi_{f1})^{0.0029} \wedge (\phi_{f2})^{0.99}$, where $\phi_{f1} = (\Diamond_{[0,531]}^{w^{r1}} \Box_{[0,531]}^{w^{r2}} \pi^{r1}) \wedge (\Diamond_{[0,531]}^{w^{s1}} \Box_{[0,531]}^{w^{s2}} \pi^{s1}) \wedge (\Diamond_{[0,531]}^{w^{d1}} \Box_{[0,531]}^{w^{d2}} \pi^{d1})$, $\phi_{f2} = (\Box_{[0,531]}^{w^r} \pi^{r2})^{0.69} \wedge (\Box_{[0,531]}^{w^s} \pi^{s2})^{0.018} \wedge (\Box_{[0,531]}^{w^d} \pi^{d2})^{0.29}$, and $\pi^{ri}, \pi^{si}, \pi^{di}, i = 1, 2$ are atomic predicates describing raw data (d), spectral data (\hat{d}), and derivative data (\dot{d}), respectively. Notice that the subformula ϕ_{f1} has an associated weight of 0.0029, which is much smaller than the 0.99 of ϕ_{f2} , indicating that ϕ_{f1} is insignificant in classifying the data. To validate this, we only use ϕ_{f2} to make prediction of the test data. The accuracy remains the same after dropping ϕ_{f1} , which also demonstrates the property of impact ordering. Hence ϕ_{f2} is utilized for interpretability analysis. The weights

of ϕ_{f2} indicate that the raw data and the derivative data are more important. As a result, $\pi^{r2} = (d \geq 0.20), \pi^{d2} = (\tilde{d} \geq 0.0067)$. If we consider the major subformulas, ϕ reads as “The distance is larger than 0.20 during 477 to 481 and the derivative is larger than 0.0067 from 228 to 229”. The top-5 weights for w^r are $[w_{477}^r, w_{478}^r, w_{479}^r, w_{480}^r, w_{481}^r] = [0.0050, 0.080, 0.82, 0.041, 0.010]$. The decision rule from π^{r2} on the raw data is shown in Fig. 6 (top). We could observe during 477 to 481, the Fly data’s distance is above 0.20, and the Beetle data’s distance is below 0.20. The top-2 weights for w^d are $[w_{228}^d, w_{229}^d] = [0.99, 0.0010]$. The decision rule from π^{d2} on the derivative data is shown in Fig. 6 (bottom). We could observe during 228 to 229, the derivative of the Beetle data is below 0.0067, and the derivative of the Fly data is above 0.0067. This demonstrates the learned patterns from NSTSC can discriminate the two classes of data.

C. Ablation Study

This subsection investigates the necessity of multi-view data representation in the TSC process. The node phase classifier and tree phase classifier are obviously necessary. To demonstrate the necessity of each module, we exclude it and show the degradation of the model performance. First, we study the effect of each data representation by keeping the other two representations. Second, we only use the raw data to evaluate NSTSC. Table V shows the results of the ablation study on the UCR dataset. Specifically, we have the following findings.

(1) Removing the raw, spectral, or derivative data degrades the average accuracy to 79.13%, 79.61%, 79.85%, respectively. This validates the multi-view data representation facilitates the learning process.

(2) The average accuracy of NSTSC with raw data only is 79.36%, which indicates using a single raw data representation is still effective for TSC.

D. Interpretability Comparison with SAX-VSM

Existing interpretable TSC models, such as Symbolic Aggregate Approximation and Vector Space Model (SAX-VSM) [35], exploit subsequences of TSD as features and find the discriminatory subsequences to represent a class. Nevertheless, it cannot provide a human-readable formula that is intuitively interpretable. By contrast, NSTSC can both identify the discriminatory subsequences and provide an interpretable formula analogous to natural language. The “GunPoint” dataset is utilized to validate this statement, which involves two actors making a motion with their hands. For the “Gun” class, the actors first move their hands above a hip-mounted holster and then move their hands down to grasp the gun and move their hands up to the shoulder level. For the “Point” class, the actors directly move their fingers up to the shoulder level. The x position of the actors’ hands is tracked. Fig. 7 shows the intervals extracted by NSTSC and SAX-VSM. It is evident that intervals from NSTSC align with the subsequences identified by SAX-VSM because of the intersection. The patterns correspond to the difference between moving down to grasp the gun and directly moving up. Furthermore, if

we truncate the small weights, NSTSC can provide a wSTL formula $\phi = (x(34) \leq -0.060) \vee (x(35) \leq 0.040) \vee (x(36) \leq 0.28) \vee (x(109) \leq -0.48) \vee (x(110) \leq -0.65)$, which reads in plain English as “ x pos at 34 s is below -0.060 or x pos at 35 s is below 0.04 or x pos at 36 s is below 0.28 or x pos at 109 s is below -0.48 or x pos at 110 s is below -0.65”, while SAX-VSM cannot provide such human-readable formulas. Note Table III excludes the comparison with SAX-VSM as the accuracy is not available in the original paper.

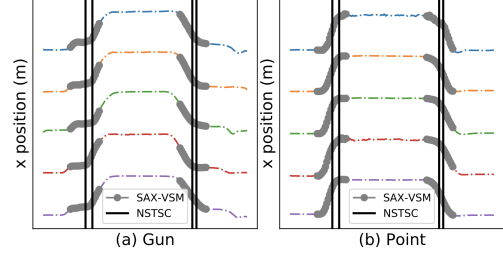


Fig. 7. Characteristic patterns in NSTSC and SAX-VSM.

VII. CONCLUSION

In this study, we propose a neuro-symbolic model, NSTSC, to classify TSD. We propose a novel weighted semantics for STL to quantify the truth degree of a wSTL formula over TSD. The NSTSC is designed by associating each neuron with a symbolic expression in a wSTL formula such that the NSTSC is differentiable and can be expressed as human-readable statements. A decision tree-based classifier is designed to learn formula structures and classify multiclass TSD. We apply NSTSC on two groups of time-series datasets to demonstrate that it can achieve comparable performance with state-of-the-art models. Furthermore, the interpretability analysis demonstrates that NSTSC can produce logical formulas that are easy to understand and match with domain knowledge, which is a property that existing models do not possess.

VIII. ACKNOWLEDGEMENTS

This research is sponsored by the Rensselaer-IBM AI Research Collaboration (<http://airc.rpi.edu>), part of the IBM AI Horizons Network; the National Science Foundation under Grant CMMI-1936578; and the Defense Advanced Research Projects Agency (DARPA) through Cooperative Agreement D20AC00004 awarded by the U.S. Department of the Interior (DOI), Interior Business Center. The content of the information does not necessarily reflect the position or the policy of the Government, and no official endorsement should be inferred.

REFERENCES

- [1] A. Dempster, D. F. Schmidt, and G. I. Webb, “Minirocket: A very fast (almost) deterministic transform for time series classification,” in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 248–257.
- [2] A. P. Ruiz, M. Flynn, J. Large, M. Middlehurst, and A. Bagnall, “The great multivariate time series classification bake off: a review and experimental evaluation of recent algorithmic advances,” *Data Mining and Knowledge Discovery*, vol. 35, no. 2, pp. 401–449, 2021.
- [3] S. Selvin, R. Vinayakumar, E. Gopalakrishnan, V. K. Menon, and K. Soman, “Stock price prediction using lstm, rnn and cnn-sliding window model,” in *2017 international conference on advances in computing, communications and informatics*. IEEE, 2017, pp. 1643–1647.

- [4] J. de Jong, M. A. Emon, P. Wu, R. Karki, M. Sood, P. Godard, A. Ahmad, H. Vrooman, M. Hofmann-Apitius, and H. Fröhlich, “Deep learning for clustering of multivariate clinical patient trajectories with missing values,” *GigaScience*, vol. 8, no. 11, p. giz134, 2019.
- [5] Q. Ma, C. Chen, S. Li, and G. W. Cottrell, “Learning representations for incomplete time series clustering,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, 2021, pp. 8837–8846.
- [6] A. Shoeb and J. Guttag, “Application of machine learning to epileptic seizure detection,” in *ICML*, 2010.
- [7] A. Subasi, J. Kevric, and M. A. Canbaz, “Epileptic seizure detection using hybrid machine learning methods,” *Neural Computing and Applications*, vol. 31, no. 1, pp. 317–325, 2019.
- [8] G. Yang, C. Belta, and R. Tron, “Continuous-time signal temporal logic planning with control barrier functions,” in *2020 American Control Conference (ACC)*. IEEE, 2020, pp. 4612–4618.
- [9] M. Charitidou and D. V. Dimarogonas, “Signal temporal logic task decomposition via convex optimization,” *IEEE Control Systems Letters*, 2021.
- [10] R. Yan, Z. Xu, and A. Julius, “Swarm signal temporal logic inference for swarm behavior analysis,” *IEEE Robotics and Automation Letters*, vol. 4, no. 3, pp. 3021–3028, 2019.
- [11] G. Bombara and C. Belta, “Offline and online learning of signal temporal logic formulae using decision trees,” *ACM Transactions on Cyber-Physical Systems*, vol. 5, no. 3, pp. 1–23, 2021.
- [12] J. Hills, J. Lines, E. Baranauskas, J. Mapp, and A. Bagnall, “Classification of time series by shapelet transformation,” *Data mining and knowledge discovery*, vol. 28, no. 4, pp. 851–881, 2014.
- [13] A. Dempster, F. Petitjean, and G. I. Webb, “Rocket: exceptionally fast and accurate time series classification using random convolutional kernels,” *Data Mining and Knowledge Discovery*, vol. 34, no. 5, pp. 1454–1495, 2020.
- [14] Z. Wang, W. Yan, and T. Oates, “Time series classification from scratch with deep neural networks: A strong baseline,” in *2017 International joint conference on neural networks*. IEEE, 2017, pp. 1578–1585.
- [15] M. Middlehurst, J. Large, M. Flynn, J. Lines, A. Bostrom, and A. Bagnall, “Hive-cote 2.0: a new meta ensemble for time series classification,” *Machine Learning*, vol. 110, no. 11, pp. 3211–3243, 2021.
- [16] M. Middlehurst, J. Large, G. Cawley, and A. Bagnall, “The temporal dictionary ensemble (tde) classifier for time series classification,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2020, pp. 660–676.
- [17] M. Middlehurst, J. Large, and A. Bagnall, “The canonical interval forest (cif) classifier for time series classification,” in *2020 IEEE international conference on big data (big data)*. IEEE, 2020, pp. 188–195.
- [18] A. Shifaz, C. Pelletier, F. Petitjean, and G. I. Webb, “Ts-chief: a scalable and accurate forest algorithm for time series classification,” *Data Mining and Knowledge Discovery*, vol. 34, no. 3, pp. 742–775, 2020.
- [19] Z. Kong, A. Jones, A. Medina Ayala, E. Aydin Gol, and C. Belta, “Temporal logic inference for classification and prediction from data,” in *International conference on Hybrid systems: computation and control*, 2014, pp. 273–282.
- [20] Z. Xu and A. A. Julius, “Census signal temporal logic inference for multiagent group behavior analysis,” *IEEE Transactions on Automation Science and Engineering*, vol. 15, no. 1, pp. 264–277, 2016.
- [21] Z. Wang, W. Zhang, N. Liu, and J. Wang, “Scalable rule-based representation learning for interpretable classification,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [22] Z. Wang, W. Zhang, L. Ning, and J. Wang, “Transparent classification with multilayer logical perceptrons and random binarization,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 04, 2020, pp. 6331–6339.
- [23] R. Yan, A. Julius, M. Chang, A. Fokoue, T. Ma, and R. Uceda-Sosa, “Stone: Signal temporal logic neural network for time series classification,” in *2021 International Conference on Data Mining Workshops (ICDMW)*. IEEE, 2021, pp. 778–787.
- [24] R. Zhang, B. Ashuri, and Y. Deng, “A novel method for forecasting time series based on fuzzy logic and visibility graph,” *Advances in Data Analysis and Classification*, vol. 11, no. 4, pp. 759–783, 2017.
- [25] R. Riegel, A. Gray, F. Luus, N. Khan, N. Makondo, I. Y. Akhalwaya, H. Qian, R. Fagin, F. Barahona, U. Sharma *et al.*, “Logical neural networks,” *arXiv preprint arXiv:2006.13155*, 2020.
- [26] P. Kapanipathi, I. Abdelaziz, S. Ravishankar, S. Roukos, A. Gray, R. Astudillo, M. Chang, C. Cornelio, S. Dana, A. Fokoue *et al.*, “Question answering over knowledge bases by leveraging semantic parsing and neuro-symbolic reasoning,” *arXiv preprint arXiv:2012.01707*, 2020.
- [27] N. Mehdipour, C.-I. Vasile, and C. Belta, “Specifying user preferences using weighted signal temporal logic,” *IEEE Control Systems Letters*, vol. 5, no. 6, pp. 2006–2011, 2021.
- [28] J. Lines, S. Taylor, and A. Bagnall, “Time series classification with hive-cote: The hierarchical vote collective of transformation-based ensembles,” *ACM Transactions on Knowledge Discovery from Data*, vol. 12, no. 5, 2018.
- [29] E. J. Keogh and M. J. Pazzani, “Derivative dynamic time warping,” in *Proceedings of the 2001 SIAM international conference on data mining*. SIAM, 2001, pp. 1–11.
- [30] G. Bombara, C.-I. Vasile, F. Penedo, H. Yasuoka, and C. Belta, “A decision tree approach to data classification using signal temporal logic,” in *Proceedings of the 19th International Conference on Hybrid Systems: Computation and Control*, 2016, pp. 1–10.
- [31] H. A. Dau, A. Bagnall, K. Kamgar, C.-C. M. Yeh, Y. Zhu, S. Gharghabi, C. A. Ratanamahatana, and E. Keogh, “The ucr time series archive,” *IEEE/CAA Journal of Automatica Sinica*, vol. 6, no. 6, pp. 1293–1305, 2019.
- [32] M. P. Caley, V. L. Martins, and E. A. O’Toole, “Metalloproteinases and wound healing,” *Advances in wound care*, vol. 4, no. 4, pp. 225–234, 2015.
- [33] P. Schäfer, “The BOSS is concerned with time series classification in the presence of noise,” *Data Mining and Knowledge Discovery*, vol. 29, no. 6, pp. 1505–1530, 2015.
- [34] N. Cabello, E. Naghizade, J. Qi, and L. Kulik, “Fast and accurate time series classification through supervised interval search,” in *2020 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2020, pp. 948–953.
- [35] P. Senin and S. Malinchik, “Sax-vsm: Interpretable time series classification using sax and vector space model,” in *IEEE international conference on data mining*. IEEE, 2013, pp. 1175–1180.

A. Conversion from a Tree Classifier to a Formula

The wSTL formula conversion procedures are presented in Algorithm 3.

Algorithm 3 Tree Classifier to Formulas - *TCTF()*

Input: A tree phase classifier \mathcal{T} learned via Algorithm 2

Output: wSTL formulas for C classes, $\Phi = \{\phi_1, \dots, \phi_C\}$

- 1: Initialize a subformula list for each data class, $\phi_c = []$
 - 2: **for** each path \mathcal{P} in \mathcal{T} **do**
 - 3: **for** each Node n in \mathcal{P} **do**
 - 4: **if** Node n is the root node **then**
 - 5: $\phi_{pt} = \emptyset$
 - 6: **else if** Node n is a left child node **then**
 - 7: Set $\phi_{pt} = \phi_{pt} \wedge \phi^{n-1,*}$
 - 8: **else if** Node n is a right child node **then**
 - 9: Set $\phi_{pt} = \phi_{pt} \wedge \neg \phi^{n-2,*}$
 - 10: **end if**
 - 11: **end for**
 - 12: Append ϕ_{pt} into ϕ_c for class c represented by the leaf node of \mathcal{P}
 - 13: **end for**
 - 14: **for** $c = 1, \dots, C$ **do**
 - 15: Set $\phi_c = \emptyset$
 - 16: **for** each subformula ϕ_c^s in ϕ_c **do**
 - 17: $\phi_c \leftarrow \phi_c \vee \phi_c^s$
 - 18: **end for**
 - 19: $\Phi \leftarrow \Phi \cup \phi_c$
 - 20: **end for**
 - 21: **return** Φ
-