



Interval Temporal Logic Decision Tree Learning

Andrea Brunello² , Guido Sciavicco¹ , and Ionel Eduard Stan²

¹ Department of Mathematics and Computer Science,
University of Ferrara, Ferrara, Italy
guido.sciavicco@unife.it

² Department of Mathematics, Computer Science, and Physics,
University of Udine, Udine, Italy
andrea.brunello@uniud.it, stan.ioneleduard@spes.uniud.it

Abstract. Decision trees are simple, yet powerful, classification models used to classify categorical and numerical data, and, despite their simplicity, they are commonly used in operations research and management, as well as in knowledge mining. From a logical point of view, a decision tree can be seen as a structured set of logical rules written in propositional logic. Since knowledge mining is rapidly evolving towards temporal knowledge mining, and since in many cases temporal information is best described by interval temporal logics, propositional logic decision trees may evolve towards interval temporal logic decision trees. In this paper, we define the problem of interval temporal logic decision tree learning, and propose a solution that generalizes classical decision tree learning.

Keywords: Decision trees · Interval temporal logics · Symbolic learning

1 Introduction

It is commonly recognized that modern decision trees are of primary importance among classification models [30]. They owe their popularity mainly to the fact that they can be trained and applied efficiently even on big datasets, and that they are easily interpretable, meaning that they are not only useful for prediction per se, but also for understanding the reasons behind the predictions. Interpretability is of extreme importance in domains in which understanding the classification process is as important as the accuracy of the classification itself, such in the case of production business systems or in the computer-aided medicine domain. A typical decision tree is constructed in a recursive manner, following the traditional *Top Down Induction of Decision Trees* (TDIDT) approach [26]: starting from the root, at each node the attribute that best partitions the training data, according to a predefined score, is chosen as a test to guide the partitioning of instances into child nodes. The process continues until a sufficiently high degree of purity (with respect to the target class), or a

minimum cardinality constraint (with respect to the number of instances reaching the node), is achieved in the generated partitions. This is the case of the well-known decision tree learning algorithm ID3 [26], which is the precursor of the commonly-used C4.5 [27]. A decision tree can be seen as a structured set of rules: every node of the tree can be thought of as a decision point, and, in this way, each branch becomes a conjunction of such conditional statements, that is, a *rule*, whose right-hand part is the class. A conditional statement may have many forms: it can be a yes/no statement (for binary categorical attributes), a categorical value statement (for non-binary categorical attributes), or a splitting value statement (for numerical attributes); the arity of the resulting tree is two if all attributes are binary or numerical, or more, if there are categorical attributes with more than two categories. Each statement can be equivalently represented with *propositional letters*, so that a decision tree can be also seen as a structured set of *propositional logic* rules.

Temporal Classification: Static Solutions. Just focusing on the static aspects of data is not always adequate for classification; for example, in the medical domain, one may want to take into account which symptoms a patient was experiencing at the same time, or whether two symptoms were overlapping. That is, in some application domains, the temporal aspects of the information may be essential to an accurate prediction. Within static decision tree learning, temporal information may be aggregated in order to circumvent the absence of explicit tools for dealing with temporal information (for example, a patient can be labelled with a natural number describing how many times he/she has been running a fever during the observation period); the ability of a decision tree to perform a precise classification based on such processed data, however, strongly depends on how well data are prepared, and therefore on how well the underlying domain is understood. Alternatively, decision trees have been proposed that use *frequent patterns* [15, 19, 22] in nodes, considering the presence/absence of a frequent pattern as a categorical attribute [13, 18]. Nevertheless, despite being the most common approach to (explicit) temporal data classification, frequent patterns in sequences or series have a limited expressive power, as they are characterized by being *existential* and by intrinsically representing temporal information with instantaneous events.

Our Approach: Interval Temporal Logic Decision Trees. A different approach to temporal classification is mining temporal logic formulas, and since temporal databases universally adopt an interval-based representation of time, the ideal choice to represent temporal information in data is *interval* temporal logic. The most representative propositional interval temporal logic is Halpern and Shoham’s Modal Logic of Allen’s Relations [20], also known as HS. Its language encompasses one modal operator for each interval-to-interval relation, such as *meets* or *before*, and the computational properties of HS and its fragments have been studied in the recent literature (see, e.g. [10–12]). The very high expressive power of HS, as well as its versatility, make HS the ideal candidate to serve as the basis of a temporal decision tree learning algorithm. Based on these premises, we propose in this paper a decision tree learning algorithm that produces HS-based

trees. Our proposal, Temporal ID3, is a direct generalization of the ID3 algorithm [26], founded on the logical interpretation of tree nodes, and focuses on data representation and node generation; we borrow other aspects, such as splitting based on information gain and the overall learning process from the original algorithm. The accuracy of a decision tree and its resilience to over-fitting also depends on the stopping criterion and possible post-pruning operations, but we do not discuss these aspects here.

Existing Approaches to Temporal Logic Decision Trees. Learning temporal logic decision trees is an emerging field in the analysis of physical systems, and, among the most influential approaches, we mention learning of automata [3] and learning Signal Temporal Logic (STL) formulas [6, 14, 24, 28]. In particular, STL is a point-based temporal logic with *until* that encompasses certain metric capabilities, and learning formulas of STL has been focused on both the fine tuning of the metric parameters of a predefined formula and the learning the innermost structure of a formula; among others, decision trees have been used to this end [8]. Compared with STL decision tree learning, our approach has the advantage of learning formulas written in a well-known, highly expressive interval-based temporal logic language; because of the nature of the underlying language and of the interval temporal logic models, certain application domains fit naturally into this approach. Moreover, since our solution generalizes the classical decision tree learning algorithm ID3, and, particularly, the notion of information gain, it is not limited to binary classification only. Moreover, in [7] a first-order framework for TDIDT is presented with the attempt to make such paradigm more attractive to inductive logic programming (ILP). Such a framework provides a sound basis for logical decision tree induction; in opposition, we employ the framework to represent *modal*, instead of first-order, relational data. Additionally, our approach should not be confused with [23], in which the term *interval* indicates an uncertain numerical value (e.g., *the patient has a fever of 38 Celsius* versus *the patient has a fever between 37.5 and 38.5 Celsius*), and in which an algorithm for inducing decision trees on such uncertain data is presented that is based on the so-called Kolmogorov-Smirnov criterion, but the data that are object of that study are not necessarily temporal, and the produced trees do not employ any temporal (logical) relation. In [4, 29] and [21], the authors present two other approaches to a temporal generalization of decision tree learning. In the former, the authors provide a general method for building point-based temporal decision trees, but with no particular emphasis on any supporting formal language. In the latter, the constructed trees can be seen as real-time algorithms that have the ability to make decisions even if the entire description of the instance is not yet available. Finally, in [16], a generalization of the decision tree model is presented in which nodes are possibly labelled with a timestamp to indicate when a certain condition should be checked.

Summarizing, our approach is essentially different from those presented in the literature in several aspects. As a matter of fact, by giving a logical perspective to decision tree learning, we effectively generalize the learning model to a temporal one, instead of introducing a new paradigm. In this way, instances that

present some temporal component are naturally seen as timelines, and, thanks to the expressive power provided by HS, our algorithm can learn a decision tree based on the temporal relations between values, instead of the static information carried by the values.

2 Preliminaries

Decision Trees. Decision tree induction is based on the following simple concepts [27]. Given a set of observable values $V = \{v_1, v_2, \dots, v_n\}$, with associated probabilities $\Pi = \{\pi_1, \pi_2, \dots, \pi_n\}$, the *information conveyed by Π* (or *entropy*), is defined as:

$$E(\Pi) = - \sum_{i=1}^n \pi_i \log(\pi_i).$$

Assume that a dataset \mathcal{T} has n instances, each characterized by the attributes A_1, \dots, A_l , and distributed over s classes C_1, \dots, C_s . Each class C can be seen as the subset of \mathcal{T} composed of precisely those instances classified as C , so that the information needed to identify the class of an element of \mathcal{T} is:

$$Info(\mathcal{T}) = E(\{\frac{|C_1|}{|\mathcal{T}|}, \frac{|C_2|}{|\mathcal{T}|}, \dots, \frac{|C_s|}{|\mathcal{T}|}\}).$$

Intuitively, the entropy is inversely proportional to the purity degree of \mathcal{T} with respect to the class values. *Splitting*, which is the main operation in decision tree learning, is performed over a specific attribute A . If A is categorical and its domain $Dom(A)$ consists of m distinct values, we can split \mathcal{T} into $\mathcal{T}_1, \dots, \mathcal{T}_m$, each \mathcal{T}_i being characterized by A having precisely the value a_i (i.e., $A = a_i$). The information of a categorical split, therefore, is:

$$InfoCat(A, \mathcal{T}) = \sum_{i=1}^m (\frac{|\mathcal{T}_i|}{|\mathcal{T}|}) Info(\mathcal{T}_i).$$

If, on the other hand, A is numerical, then the set $\{a_1 < \dots < a_m\}$ of *actual* values for A that are present in \mathcal{T} gives rise to $m - 1$ possible splits, all of them binary, and the information conveyed by each possible split is, then, a function not only of the attribute but also of the chosen value:

$$InfoNum(A, a_i, \mathcal{T}) = (\frac{|\mathcal{T}_1|}{|\mathcal{T}|}) Info(\mathcal{T}_1) + (\frac{|\mathcal{T}_2|}{|\mathcal{T}|}) Info(\mathcal{T}_2),$$

where \mathcal{T}_1 (respectively, \mathcal{T}_2) encompasses all and only those instances with $A \leq a_i$ (respectively, $A > a_i$). The information conveyed by an attribute can be consequently defined as:

HS	Allen's relations	Graphical representation
$\langle A \rangle$	$[x, y]R_A[x', y'] \Leftrightarrow y = x'$	
$\langle L \rangle$	$[x, y]R_L[x', y'] \Leftrightarrow y < x'$	
$\langle B \rangle$	$[x, y]R_B[x', y'] \Leftrightarrow x = x', y' < y$	
$\langle E \rangle$	$[x, y]R_E[x', y'] \Leftrightarrow y = y', x < x'$	
$\langle D \rangle$	$[x, y]R_D[x', y'] \Leftrightarrow x < x', y' < y$	
$\langle O \rangle$	$[x, y]R_O[x', y'] \Leftrightarrow x < x' < y < y'$	

Fig. 1. Allen's interval relations and HS modalities.

$$InfoAtt(A, T) = \begin{cases} InfoCat(A, T) & \text{if } A \text{ is categorical,} \\ \min_{a_i \in Dom(A)} \{InfoNum(A, a_i, T)\} & \text{if } A \text{ is numerical,} \end{cases}$$

and the *information gain* brought by A is defined as:

$$Gain(A, T) = Info(T) - InfoAtt(A, T).$$

The information gain, which can be also seen as the reduction of the expected entropy when the attribute A has been chosen, is used to drive the splitting process, that is, to decide over which attribute (and how) the next split is performed. The underlying principle to decision tree building consists of recursively splitting the dataset over the attribute that guarantees the greatest information gain until a certain stopping criterion is met. Each split can be seen as a propositional condition *if p then -, else -*. When splitting is performed over a numerical attribute, e.g., $A \leq a_i$, then the corresponding propositional statement is simply the condition itself (in our example, is a propositional letter $p_{A \leq a_i}$); when it is performed over a categorical attribute, e.g., $A = a_1, A = a_2, \dots$, then each statement is a propositional statement (in our example, $p_{A=a_1}, p_{A=a_2}, \dots$) on its own.

Interval Temporal Logic. Let $\mathbb{D} \subseteq \mathbb{N}$. In the *strict* interpretation, an *interval* over \mathbb{D} is an ordered pair $[x, y]$, where $x, y \in \mathbb{D}$ and $x < y$, and we denote by $\mathbb{I}(\mathbb{D})$ the set of all intervals over \mathbb{D} . If we exclude the identity relation, there are 12 different Allen's relations between two intervals in a linear order [1]: the six relations R_A (adjacent to), R_L (later than), R_B (begins), R_E (ends), R_D (during), and R_O (overlaps), depicted in Fig. 1, and their inverses, that is, $R_{\bar{X}} = (R_X)^{-1}$, for each $X \in \mathcal{A}$, where $\mathcal{A} = \{A, L, B, E, D, O\}$. Halpern and Shoham's modal logic of temporal intervals (HS) is defined from a set of propositional letters \mathcal{AP} , and by associating a universal modality $[X]$ and an existential one $\langle X \rangle$ to each Allen's relation R_X . Formulas of HS are obtained by

$$\varphi ::= p \mid \neg\varphi \mid \varphi \vee \varphi \mid \langle X \rangle\varphi \mid \langle \bar{X} \rangle\varphi,$$

where $p \in \mathcal{AP}$ and $X \in \mathcal{A}$. The other Boolean connectives and the logical constants, e.g., \rightarrow and \top , as well as the universal modalities $[X]$, can be defined in the standard way. For each $X \in \mathcal{A}$, the modality $\langle \bar{X} \rangle$ (corresponding to the inverse relation $R_{\bar{X}}$ of R_X) is said to be the *transpose* of the modalities $\langle X \rangle$, and vice versa. The semantics of HS formulas is given in terms of *timelines* $T = \langle \mathbb{I}(\mathbb{D}), V \rangle$ ¹, where \mathbb{D} is a linear order and $V : \mathcal{AP} \rightarrow 2^{\mathbb{I}(\mathbb{D})}$ is a *valuation function* which assigns to each atomic proposition $p \in \mathcal{AP}$ the set of intervals $V(p)$ on which p holds. The *truth* of a formula φ on a given interval $[x, y]$ in an interval model T is defined by structural induction on formulas as follows:

$$\begin{aligned} T, [x, y] \models p & \quad \text{if } [x, y] \in V(p), \text{ for } p \in \mathcal{AP}; \\ T, [x, y] \models \neg\psi & \quad \text{if } T, [x, y] \not\models \psi; \\ T, [x, y] \models \psi \vee \xi & \quad \text{if } T, [x, y] \models \psi \text{ or } T, [x, y] \models \xi; \\ T, [x, y] \models \langle X \rangle\psi & \quad \text{if there is } [w, z] \text{ s.t. } [x, y]R_X[w, z] \text{ and } T, [w, z] \models \psi; \\ T, [x, y] \models \langle \bar{X} \rangle\psi & \quad \text{if there is } [w, z] \text{ s.t. } [x, y]R_{\bar{X}}[w, z] \text{ and } T, [w, z] \models \psi. \end{aligned}$$

HS is a very general interval temporal language and its satisfiability problem is undecidable [20]. Our purpose here, however, is to study the problem of formula *induction* in the form of decision trees, and not of formula *deduction*, and therefore the computational properties of the satisfiability problem can be ignored at this stage.

3 Motivations

In this section, we present some realistic scenarios in which learning a temporal decision tree may be convenient, and, then, we discuss aspects of data preprocessing related to the temporal component.

Learning. There are several application domains in which learning a temporal decision tree may be useful. Consider, for example, a medical scenario in which we consider a dataset of classified patients, each one characterized by its medical history, as in Fig. 2, top. Assume, first, that we are interested in learning a *static* (propositional) classification model. The history of our patients, that is, the collection of all relevant pieces of information about tests, results, symptoms, and hospitalizations of the patient that occurred during the entire observation period, must be processed so that temporal information is subsumed in propositional letters. For instance, if some patient has been running a fever during the observation period, we may use a proposition *fever*, with positive values for those patient that have had fever, and negative values for the others (as in Fig. 2, bottom, left). Depending on the specific case, we may, instead, use the actual temperature of each patient, and a static decision tree learning system may split over *fever* $< t$, for some threshold temperature t , effectively introducing a new propositional letter, and therefore a binary split. Either way, the temporal information is lost in the preprocessing. For example, we can no longer

¹ We deliberately use the symbol T to indicate both a timeline and an instance in a dataset.

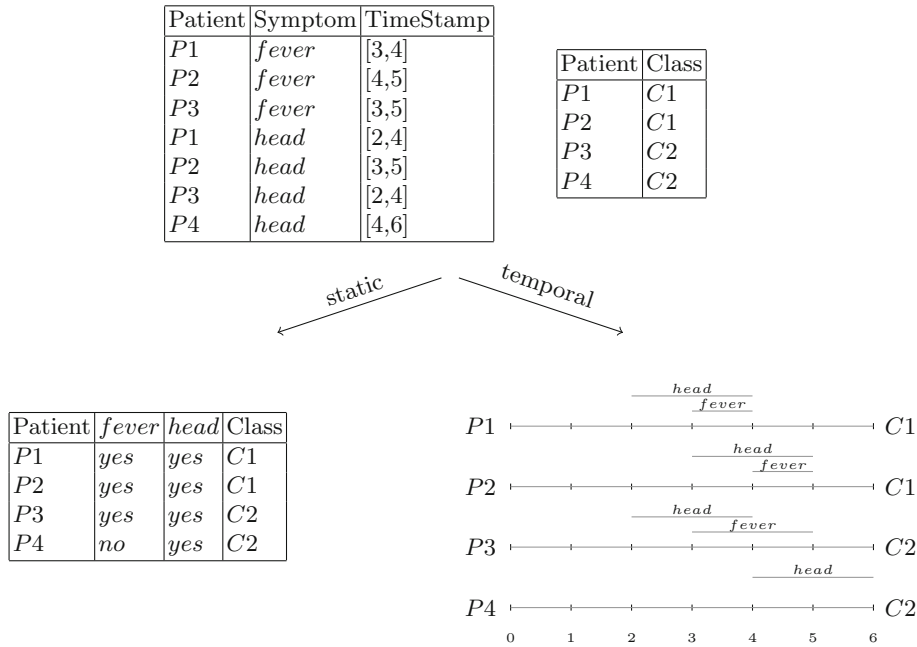


Fig. 2. Example of static and temporal treatment of information in the medical domain.

take into account whether *fever* occurred before, after, or while the patient was experiencing headache (*head*), which may be a relevant information for a classification model. By generating, instead, the timeline of each patient (as in Fig. 2, bottom, right), we keep all events and their relative qualitative relations. By learning a decision tree on a preprocessed dataset such as the one in Fig. 2 (bottom, left), we see that the attribute *head* has zero variance, and therefore zero predictive capabilities; then, we are forced to build a decision tree using attribute *fever* alone, which results in a classifier with 75% accuracy. On the contrary, by using the temporal information in the learning process, we are able to distinguish the two classes: *C1* is characterized by presenting both *head* and *fever*, but not overlapping, and this classifier has, in this toy example, 100% accuracy. In this example, the term *accuracy* refers to the *training set accuracy* (we do not consider independent training and test data), that is, the ability of the classification system to discern among classes on the data used to train the system itself; it should not be confused with *test set accuracy*, which measures the real classification performances that can be expected on future, real-life examples.

Alternatively, consider a problem in the natural language processing domain. In this scenario, a timeline may represent a *conversation* between two individuals. It is known that, in automatic processing of conversations, it is sometimes interesting to label each interval of time with one or more *context*, that is, a particular topic that is being discussed [2, 5, 25], in order to discover the exis-

tence of unexpected or interesting temporal relations among them. Suppose, for example, that a certain company wants to analyze conversations between selling agents and potential customers: the agents contact the customers with the aim of selling a certain product, and it is known that certain contexts, such as the price of the product (*price*), its known advantages (*advantages*) over other products, and its possible minor defects (*disadvantages*) are interesting. Assume that each conversation has been previously classified between those that have been successful and those that ended without the product being acquired. Now, we want to learn a model able to predict such an outcome. By using only static information, nearly every conversation would be labelled with the three contexts, effectively hiding the underlying knowledge, if it exists. By keeping the relative temporal relations between contexts, instead, we may learn useful information, such as, for example, *if price and disadvantages are not discussed together, the conversation will be likely successful*.

Preprocessing. Observe, now, how switching from static to temporal information influences data preparation. First, in a context such as the one described in our first example, numerical attributes may become less interesting: for instance, the information on *how many times* a certain symptom occurred, or its *frequency*, are not needed anymore, considering that each occurrence is taken into account in the timeline. Moreover, since the focus is on attributes *relative temporal positions*, even categorical attributes may be ignored in some contexts: for instance, in our scenario, we may be interested in establishing the predictive value of the relative temporal position of *fever* and *head* regardless of the sex or age of the patient. It is also worth underlying that propositional attributes over intervals allow us to express a variety of situations, and sometimes propositional labelling may result in *gaining* information, instead of loosing it. Consider, again, the case of fever, and suppose that a certain patient is experiencing low fever in an interval $[x, y]$, say, a given day, and that during just one hour of that day, that is, over the interval $[w, z]$ strictly contained in $[x, y]$, he/she has an episode of high fever. A natural choice is to represent such a situation by labelling the interval $[x, y]$ with *lo* and its sub-interval $[w, z]$ with *hi*. On the other hand, representing the same pieces of information as three intervals $[x, w]$, $[w, z]$, $[z, y]$ respectively labelled with *lo*, *hi*, and *lo*, which would be the case with a point-based representation (or with an interval-based representation under the homogeneity assumption), would be unnatural, and it would entail hiding a potentially important information such as: *“the patient presented low fever during the entire day, except for a brief episode of high fever”*. Building on such considerations, our approach in the rest of this paper is based on propositional, non-numerical attributes only.

4 Learning Interval Temporal Logic Decision Trees

In this section we describe a generalization of the algorithm ID3 that is capable of learning a binary decision tree over a temporal dataset, as in the examples of the previous section; as in classical decision trees, every branch of a temporal decision tree can be read as a logical formula, but instead of classical propositional logic

we use the temporal logic HS. To this end, we generalize the notion of information gain, while, at this stage, we do not discuss pre-pruning, post-pruning, and purity degree of a sub-tree [9, 27].

Data Preparation and Presentation. We assume that the input dataset contains timelines as instances. For the sake of simplicity, we also assume that all timelines are based on the same finite domain \mathbb{D} of length N (from 0 to $N-1$). The dataset \mathcal{T} can be seen as an array of n structures; $\mathcal{T}[j]$ represents the j -th timeline of the dataset, and it can be thought of as an interval model. Given a dataset \mathcal{T} , we denote by \mathcal{AP} the set of all propositional letters that occur in it.

Temporal Information. We are going to design the learning process based on the same principles of classical decision tree learning. This means that we need to define a notion of splitting as well as a notion of information conveyed by a split, and, to this end, we shall use the truth relation as defined in Sect. 2 applied to a timeline. Unlike the atemporal case, splits are not performed over attributes, but, instead, over propositional letters. Splitting is defined *relatively to an interval* $[x, y]$, and it can be local, if it is applied on $[x, y]$ itself, or temporal, in which case it depends on the existence of an interval $[z, t]$ related to $[x, y]$ and the particular relation R_X such that $[x, y]R_X[z, t]$ (or the other way around). A *local split* of \mathcal{T} into \mathcal{T}_1 and \mathcal{T}_2 , where $[x, y]$ is the *reference* interval of \mathcal{T} , and p is the propositional letter over which the split takes place is defined by:

$$\begin{aligned}\mathcal{T}_1 &= \{T \in \mathcal{T} \mid T, [x, y] \models p\}, \\ \mathcal{T}_2 &= \{T \in \mathcal{T} \mid T, [x, y] \models \neg p\}.\end{aligned}\tag{1}$$

On the contrary, a *temporal split*, in the same situation, over the temporal relation R_X , is defined by:

$$\begin{aligned}\mathcal{T}_1 &= \{T \in \mathcal{T} \mid T, [x, y] \models \langle X \rangle p\}, \\ \mathcal{T}_2 &= \{T \in \mathcal{T} \mid T, [x, y] \models [X] \neg p\}.\end{aligned}\tag{2}$$

Consequently, the *local information gain* of a propositional letter p is defined as:

$$LocalGain(p, \mathcal{T}) = Info(\mathcal{T}) - \left(\left(\frac{|\mathcal{T}_1|}{|\mathcal{T}|} \right) Info(\mathcal{T}_1) + \left(\frac{|\mathcal{T}_2|}{|\mathcal{T}|} \right) Info(\mathcal{T}_2) \right),$$

where \mathcal{T}_1 and \mathcal{T}_2 are defined as in (1), while the *temporal information gain* of a propositional letter p is defined as:

$$TemporalGain(p, \mathcal{T}) = Info(\mathcal{T}) - \min_{X \in \mathcal{A}} \left\{ \left(\frac{|\mathcal{T}_1|}{|\mathcal{T}|} \right) Info(\mathcal{T}_1) + \left(\frac{|\mathcal{T}_2|}{|\mathcal{T}|} \right) Info(\mathcal{T}_2) \right\},$$

where \mathcal{T}_1 and \mathcal{T}_2 are defined as in (2) and depend on the relation R_X . Therefore, the information gain of a propositional letter becomes:

$$Gain(p, \mathcal{T}) = \max\{LocalGain(p, \mathcal{T}), TemporalGain(p, \mathcal{T})\},$$

and, at each step, we aim to find the letter that maximizes the gain.

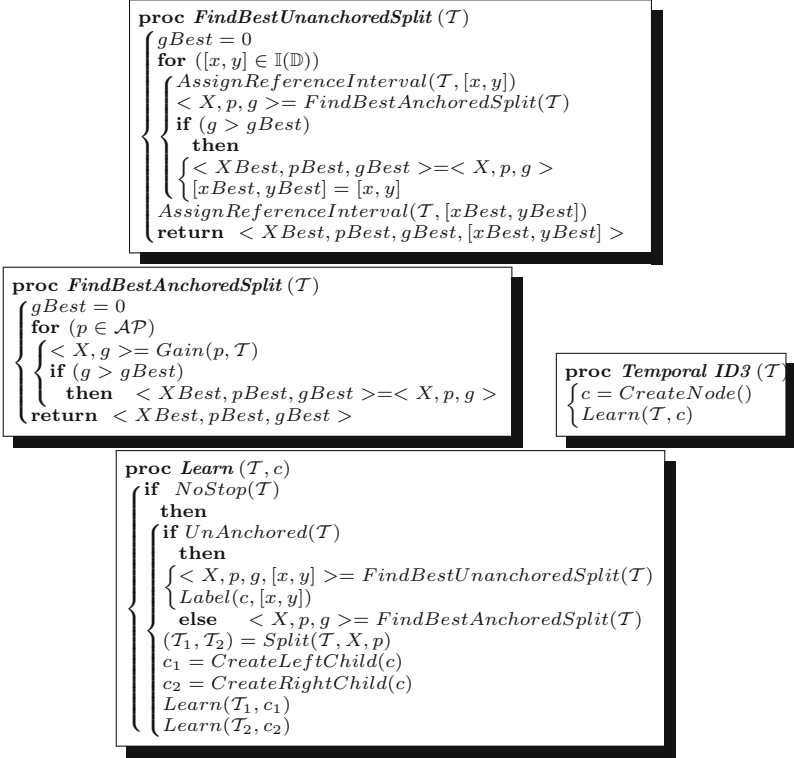


Fig. 3. The algorithm Temporal ID3.

The Algorithm. Let us analyze the code in Fig. 3. At the beginning, the timelines in T are not assigned any reference interval, and we say that the dataset is *unanchored*. The procedure *FindBestUnanchoredSplit* systematically explores every possible reference interval of an unanchored dataset, and, for each one of them, calls *FindBestAnchoredSplit*, which, in turn, tries every propositional letter (and, implicitly, every temporal relation) in the search of the best split. This procedure returns the best possible triple $\langle X, p, g \rangle$, where X is an interval relation, if the best split is temporal, or it has no value, if the best split is local, p is a propositional letter, and g is the information gain. *Temporal ID3* first creates a root node, and then calls *Learn*. The latter, in turn, first checks possible stopping conditions, and then finds the best split into two datasets T_1 and T_2 . Of these, the former is now *anchored* (to the reference interval returned by *FindBestUnanchoredSplit*), while the latter is still *unanchored*. During a recursive call, when T_1 is analyzed to find its best split, the procedure for this task will be *FindBestAnchoredSplit*, called directly, instead of passing through *FindBestUnanchoredSplit*. So, in our learning model, all splits are binary. Given a node, the ‘lefthand’ outgoing edge is labeled with the chosen $\langle X \rangle p$, or just p , when the

Patient	Symptom	TimeStamp
P1	fever	[0,2]
P2	fever	[1,3]
P3	head	[0,2]
P4	head	[1,3]

Patient	Class
P1	C1
P2	C1
P3	C2
P4	C2

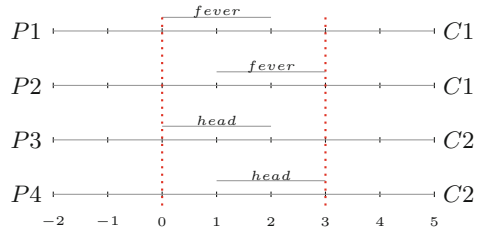


Fig. 4. Example of a problematic dataset.

split is local, whereas the corresponding ‘righthand’ edge is labeled with $[X]\neg p$ (or just $\neg p$); also, the node is labeled with a new reference interval if its corresponding dataset is unanchored. After a split, every $T \in \mathcal{T}_1$ (the existential dataset, which is now certainly anchored) is associated with a new *witnessing* interval: in fact, those instances satisfy $\langle X \rangle p$ on $[x, y]$, and, for each one of them, there is a possibly distinct witness. Witnesses are assigned by the function *Split*; while the witnessing interval of an instance may change during the process, its reference interval is set only once.

Consider, now, the function *AssignReferenceInterval* and the example shown in Fig. 4. As can be seen, neglecting the temporal dimension, one may classify the instances with just a single split based on the presence of the symptom fever (or headache). On the contrary, given the temporal dataset with domain $\mathbb{D} = \{0, 1, 2, 3\}$ it is not possible discriminate the classes within a single step. A natural solution consists of augmenting \mathbb{D} in such a way to simulate the behaviour of an infinite domain model. In our example, it suffices to consider $\mathbb{D} = \{-2, -1, 0, 1, 2, 3, 4, 5\}$, so that a single split may be based on the rule: $\langle L \rangle fever \rightarrow C1$, otherwise $C2$ holding on $[-2, -1]$ (or, equivalently, its inverse formulation on $[4, 5]$). Thus, the function *AssignReferenceIntervals*, while searching all possible reference intervals, takes into consideration two extra points at each side of the domain. Although it is possible to obtain a similar result by adding less than four points (in our example, -2 and -1 suffice), this is no longer true if we include the possibility that Temporal ID3 is called on a *subset* of HS modalities, for example, for computational efficiency reasons. Adding four points, on the other hand, guarantees that the most discriminative split can always be found.

Analysis. We now analyze the computational complexity of Temporal ID3. To this end, we first compute the cost of finding the best splitting. Since the cardinality of the domain of each timeline is N , there are $O(N^2)$ possible intervals. This means that, fixed a propositional letter and a relation R_X , computing \mathcal{T}_1 and \mathcal{T}_2 costs $O(nN^2)$, where n is the number of timelines. Therefore, the cost of *FindBestAnchoredSplit* is obtained by multiplying the cost of a single

(tentative) splitting by the number of propositional letters and the number of temporal relations (plus one, to take into account the local splitting), which sums up to $O(13nN^2|\mathcal{AP}|)$. The cost of *FindBestUnanchoredSplit* increases by a factor of N^2 , as the **for** cycle ranges over all possible intervals, and therefore it becomes $O(13nN^4|\mathcal{AP}|)$. We can increase the efficiency of the implementation by suitably pre-compute the value of $\langle X \rangle p$ for each temporal relation, each propositional letter, and each interval, thus eliminating a factor of N^2 from both costs.

If we consider \mathcal{AP} as fixed, and N as a constant, the cost of finding the best splitting becomes $O(n)$, and, under such (reasonable) assumption, we can analyze the complexity of an execution of *Learn* in terms of the number n of timelines. Two cases are particularly interesting. In the worst case, every binary split leads to a very unbalanced partition of the dataset, with $|\mathcal{T}_1| = 1$ and $|\mathcal{T}_2| = n - 1$ (or the other way around). The recurrence that describes such a situation is:

$$t(n) = t(n - 1) + O(n),$$

which can be immediately solved to obtain $t(n) = O(n^2)$. However, computing the worst case has only a theoretical value; we can reasonably expect Temporal ID3 to behave like a randomized divide-and-conquer algorithm, and its computational complexity to tend towards the average case. In the average case, every binary split leads to a non-unbalanced partition, but we cannot foresee the relative cardinality of each side of the partition. Assuming that every partition is equally probable, the recurrence that describes this situation is:

$$t(n) = \frac{1}{n} \sum_{k=1}^n (t(k) + t(n - k)) + O(n).$$

We want to prove that $t(n) = O(n \log(n))$. To this end, we first prove a useful bound for the expression $\sum_{k=1}^n k \log(k)$, as follows:

$$\begin{aligned} \sum_{k=1}^n (k \log(k)) &= \sum_{k=1}^{\lceil \frac{n}{2} \rceil - 1} (k \log(k)) + \sum_{k=\lceil \frac{n}{2} \rceil}^n (k \log(k)) \\ &\leq \sum_{k=1}^{\lceil \frac{n}{2} \rceil - 1} (k \log(\frac{n}{2})) + \sum_{k=\lceil \frac{n}{2} \rceil}^n (k \log(n)) \\ &= (\log(n) - 1) \sum_{k=1}^{\lceil \frac{n}{2} \rceil - 1} k + \log(n) \sum_{k=\lceil \frac{n}{2} \rceil}^n k \\ &= \log(n) \sum_{k=1}^n k - \sum_{k=1}^{\lceil \frac{n}{2} \rceil - 1} k \\ &\leq \frac{1}{2} \log(n) n(n + 1) - \frac{1}{2} \frac{n}{2} \left(\frac{n}{2} + 1 \right) \\ &= \frac{1}{2} (n^2 \log(n) + n \log(n)) - \frac{1}{8} n^2 - \frac{1}{4} n. \end{aligned}$$

Now, we prove, by induction, that $t(n) \leq an \log(n) + b$ for some positive constants a, b , as follows:

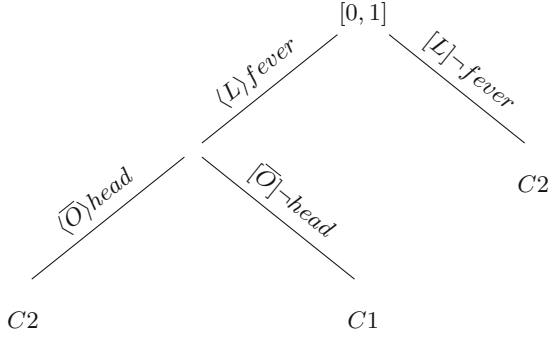


Fig. 5. A decision tree learned by Temporal ID3 on the example in Fig. 2.

$$\begin{aligned}
 t(n) &= \frac{1}{n} \sum_{k=1}^n (t(k) + t(n-k)) + O(n) \\
 &= \frac{2}{n} \sum_{k=1}^n t(k) + O(n) \\
 &\leq \frac{2}{n} \sum_{k=1}^n (ak \log(k) + b) + O(n) && \text{inductive hypothesis} \\
 &= \frac{2}{n} \sum_{k=1}^n (ak \log(k)) + \frac{2}{n} \sum_{k=1}^n b + O(n) \\
 &= \frac{2a}{n} \sum_{k=1}^n (k \log(k)) + 2b + O(n) \\
 &\leq \frac{2a}{n} \left(\frac{1}{2} (n^2 \log(n) + n \log(n)) - \frac{1}{8} n^2 - \frac{1}{4} n \right) \\
 &\quad + 2b + O(n) && \text{proved above} \\
 &= an \log(n) + 2a \log(n) - \frac{an}{4} - \frac{a}{2} \\
 &\quad + 2b + O(n) \\
 &\leq an \log(n) + b. && \text{if } \frac{an}{4} \geq 2a \log(n) - \frac{a}{2} + b + O(n).
 \end{aligned}$$

Example of Execution. Consider our initial example of Fig. 2, with four timelines distributed over two classes. Since this is a toy example, there are many different combination of intervals, relations, and propositional letters that give the same information gain. Figure 5 gives one possible outcome, which seems to indicate that, looking at the entire history, the class C2 is characterized by presenting headache and overlapping fever, or no fever at all.

There are several running parameters that can be modulated for an execution of Temporal ID3, and further analysis is required to understand how they influence the final result, and, particularly, the properties of the resulting classifier. The most important ones are: (i) how to behave in case of two splits with the same information gain; (ii) how to behave in case of more than one possible witness interval for a given timeline; (iii) how to behave in case of more than

one optimal reference interval for a given unanchored temporal dataset. If we allow, in all such cases, a random choice, the resulting learning algorithm is not deterministic anymore, and different executions may result in different decision trees. This is a typical situation in machine learning (e.g., in algorithms such as *k-means clustering*, or *random forest*), that involves some experience in order to meaningfully assess the results.

5 Conclusions

Classical decision trees, which are a popular class of learning algorithms, are designed to interpret categorical and numerical attributes. In decision trees, every node can be seen as a propositional letter; therefore, a decision tree can be seen as a structured set of propositional logic rules, the right-hand part of which is a class. Since classifying based on the static aspects of data is not always adequate, and since decision tree learning cannot deal with temporal knowledge in an explicit manner, we considered the problem of learning a classification model capable to combine propositional knowledge with qualitative temporal information. Towards its solution, we showed how temporal data can be prepared in a optimal way for a temporal decision tree to be learned and presented a generalization of the classical decision tree learning algorithm ID3 that is able to split the dataset based on temporal, instead of static, information, using the well-known temporal logic HS. Future work include testing our method on real data, improving the capabilities of Temporal ID3 by enriching the underlying language, and studying the effect of different pruning and stopping conditions. Moreover, it would be interesting to study adapting ID3 to other logical languages, although this may require re-designing some key elements, such as the representation of temporal datasets, or the process that underlies the splitting algorithm.

Machine learning is generically focused on a non-logical approach to knowledge representation. However, when learning should take into account temporal aspects of data, a logical approach can be associated to classical methods, and besides decision tree learning, interval temporal logics has been already proposed as a possible tool, for example, for temporal rules extraction [17]. Focusing these approaches on fragments of interval temporal logics whose satisfiability problem is decidable (and tractable) may result into an integrated systems that pairs induction and deduction of formulas, intelligent elimination of redundant rules, and automatic verification of inducted knowledge against formal requirement. Also, using a logical approach in learning may require non-standard semantics for logical formulas (e.g., fuzzy semantics, or multi-valued propositional semantics); these, in turn, pose original and interesting questions on the theoretical side concerning the computational properties of the problems associated with these logics (i.e., satisfiability), generating, *de facto*, a cross-feeding effect on the two fields.

References

1. Allen, J.F.: Maintaining knowledge about temporal intervals. *Commun. ACM* **26**(11), 832–843 (1983)
2. Alluhaibi, R.: Simple interval temporal logic for natural language assertion descriptions. In: *Proceedings of the 11th International Conference on Computational Semantics (IWCS)*, pp. 283–293 (2015)
3. Angluin, D.: Learning regular sets from queries and counterexamples. *Inf. Comput.* **75**(2), 87–106 (1987)
4. Antipov, S., Fomina, M.: A method for compiling general concepts with the use of temporal decision trees. *Sci. Tech. Inf. Process.* **38**(6), 409–419 (2011)
5. Baeza-Yates, R.: Challenges in the interaction of information retrieval and natural language processing. In: Gelbukh, A. (ed.) *CICLing 2004*. LNCS, vol. 2945, pp. 445–456. Springer, Heidelberg (2004). https://doi.org/10.1007/978-3-540-24630-5_55
6. Bartocci, E., Bortolussi, L., Sanguinetti, G.: Data-driven statistical learning of temporal logic properties. In: Legay, A., Bozga, M. (eds.) *FORMATS 2014*. LNCS, vol. 8711, pp. 23–37. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-10512-3_3
7. Blockeel, H., De Raedt, L.: Top-down induction of first-order logical decision trees. *Artif. Intell.* **101**(1–2), 285–297 (1998)
8. Bombara, G., Vasile, C., Penedo, F., Yasuoka, H., Belta, C.: A decision tree approach to data classification using signal temporal logic. In: *Proceedings of the 19th International Conference on Hybrid Systems: Computation and Control*, pp. 1–10 (2016)
9. Breiman, L., Friedman, J., Olshen, R., Stone, C.: *Classification and Regression Trees*. Wadsworth and Brooks, Monterey (1984)
10. Bresolin, D., Della Monica, D., Goranko, V., Montanari, A., Sciavicco, G.: Metric propositional neighborhood logics on natural numbers. *Softw. Syst. Model.* **12**(2), 245–264 (2013)
11. Bresolin, D., Monica, D.D., Montanari, A., Sala, P., Sciavicco, G.: Interval temporal logics over strongly discrete linear orders: expressiveness and complexity. *Theor. Comput. Sci.* **560**, 269–291 (2014)
12. Bresolin, D., Sala, P., Sciavicco, G.: On begins, meets, and before. *Int. J. Found. Comput. Sci.* **23**(3), 559–583 (2012)
13. Brunello, A., Marzano, E., Montanari, A., Sciavicco, G.: J48S: a sequence classification approach to text analysis based on decision trees. In: *Proceedings of the 24th International Conference on Information and Software Technologies (ICIST)*, pp. 240–256 (2018)
14. Bufo, S., Bartocci, E., Sanguinetti, G., Borelli, M., Lucangelo, U., Bortolussi, L.: Temporal logic based monitoring of assisted ventilation in intensive care patients. In: Margaria, T., Steffen, B. (eds.) *ISoLA 2014*. LNCS, vol. 8803, pp. 391–403. Springer, Heidelberg (2014). https://doi.org/10.1007/978-3-662-45231-8_30
15. Cheng, H., Yan, X., Han, J., Hsu, C.: Discriminative frequent pattern analysis for effective classification. In: *Proceedings of the 23rd International Conference on Data Engineering (ICDE)*, pp. 716–725 (2007)
16. Console, L., Picardi, C., Dupré, D.: Temporal decision trees: model-based diagnosis of dynamic systems on-board. *J. Artif. Intell. Res.* **19**, 469–512 (2003)

17. Della Monica, D., de Frutos-Escrig, D., Montanari, A., Murano, A., Sciavicco, G.: Evaluation of temporal datasets via interval temporal logic model checking. In: Proceedings of the 24th International Symposium on Temporal Representation and Reasoning (TIME), pp. 11:1–11:18 (2017)
18. Fan, W., et al.: Direct mining of discriminative and essential frequent patterns via model-based search tree. In: Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 230–238 (2008)
19. Fournier-Viger, P., Gomariz, A., Šebek, M., Hlosta, M.: VGEN: fast vertical mining of sequential generator patterns. In: Bellatreche, L., Mohania, M.K. (eds.) DaWaK 2014. LNCS, vol. 8646, pp. 476–488. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-10160-6_42
20. Halpern, J., Shoham, Y.: A propositional modal logic of time intervals. *J. ACM* **38**(4), 935–962 (1991)
21. Karimi, K., Hamilton, H.J.: Temporal rules and temporal decision trees: a C4.5 approach. Technical report, CS-2001-02, Department of Computer Science, University of Regina (2001)
22. Lin, W., Orgun, M.A.: Temporal data mining using hidden periodicity analysis. In: Raś, Z.W., Ohsuga, S. (eds.) ISMIS 2000. LNCS (LNAI), vol. 1932, pp. 49–58. Springer, Heidelberg (2000). https://doi.org/10.1007/3-540-39963-1_6
23. Mballo, C., Diday, E.: Decision trees on interval valued variables. *Symbolic Data Anal.* **3**(1), 8–18 (2005)
24. Nguyen, L., Kapinski, J., Jin, X., Deshmukh, J., Butts, K., Johnson, T.: Abnormal data classification using time-frequency temporal logic. In: Proceedings of the 20th International Conference on Hybrid Systems: Computation and Control, pp. 237–242 (2017)
25. Pratt-Hartmann, I.: Temporal prepositions and their logic. *Artif. Intell.* **166**(1–2), 1–36 (2005)
26. Quinlan, J.: Induction of decision trees. *Mach. Learn.* **1**, 81–106 (1986)
27. Quinlan, J.: Simplifying decision trees. *Int. J. Hum.-Comput. Stud.* **51**(2), 497–510 (1999)
28. Rajan, A.: Automated requirements-based test case generation. *SIGSOFT Softw. Eng. Notes* **31**(6), 1–2 (2006)
29. Vagin, V., Morosin, O., Fomina, M., Antipov, S.: Temporal decision trees in diagnostics systems. In: 2018 International Conference on Advances in Big Data, Computing and Data Communication Systems (icABCD), pp. 1–10 (2018)
30. Witten, I., Frank, E., Hall, M., Pal, C.: Data Mining: Practical Machine Learning Tools and Techniques. Morgan Kaufmann, Burlington (2016)