

Complete decision tree induction functionality in scikit-learn

Ir. Sven Van Hove

Thesis submitted for the degree of
Master of Science in Artificial
Intelligence, option Engineering and
Computer Science

Thesis supervisors:

Prof. dr. Jesse Davis
Prof. dr. ir. Hendrik Blockeel

Assessor:

Dr. ir. Marc Claesen

Mentor:

Elia Van Wolputte

© Copyright KU Leuven

Without written permission of the thesis supervisors and the author it is forbidden to reproduce or adapt in any form or by any means any part of this publication. Requests for obtaining the right to reproduce or utilize parts of this publication should be addressed to the Departement Computerwetenschappen, Celestijnenlaan 200A bus 2402, B-3001 Heverlee, +32-16-327700 or by email info@cs.kuleuven.be.

A written permission of the thesis supervisors is also required to use the methods, products, schematics and programmes described in this work for industrial or commercial use, and for submitting this publication in scientific contests.

Preface

I would like to thank Elia for the insightful discussions after business hours. I would also like to thank prof. Davis and prof. Blockeel for their helpful suggestions and the full jury for taking to time to read this document. My sincere gratitude also goes to my family and friends for their continued support. Finally, I would like to show appreciation towards my employer for the flexibility in my work schedule that made obtaining this degree possible.

Ir. Sven Van Hove

Contents

Preface	i
Abstract	iv
1 Introduction	1
1.1 Context	1
1.2 Goal	1
1.3 Motivation	2
1.4 Thesis structure	3
2 Literature review	5
2.1 Prerequisites	5
2.2 Scope	5
2.3 Terminology	6
2.4 A generic TDIDT algorithm	7
2.5 Conclusion	13
3 Existing implementations	15
3.1 Capabilities	15
3.2 Next steps	19
3.3 Conclusion	22
4 Implementation extensions	23
4.1 New pruning functionality	23
4.2 Supporting tools	27
4.3 Challenges and limitations	28
4.4 Result	29
5 Methodology	31
5.1 Pruning implementation validation	31
5.2 Evaluate the need for categorical attribute support	34
5.3 Conclusion	35
6 Results and discussion	37
6.1 Pruning	37
6.2 Categorical attributes	42
6.3 Binary versus non-binary trees	42
6.4 Conclusion	42
7 Conclusion	43

7.1 Contributions	43
7.2 Retrospective	43
7.3 Future work	43
A The First Appendix	47
Bibliography	49

Abstract

The `abstract` environment contains a more extensive overview of the work. But it should be limited to one page. [\[EHWP16\]](#)

Chapter 1

Introduction

1.1 Context

Decision tree induction is one of the most well-known tools in the machine learning community. Most of the theoretical groundwork was laid in the last three decades of the previous century. Researchers Leo Breiman and Ross Quinlan have been particularly influential in this space. Some well-known algorithms include the Concept Learning System (CLS) [HMS66] and ID3 [Qui79, Qui83, Qui86] by Quinlan, and Classification And Regression Trees (CART) [BFSO84] by Breiman.

Contemporary AI researchers focus a lot of their attention on neural networks and in particular deep learning — the recent hype around DeepMind’s AlphaGo [SSS⁺17] victories comes to mind. Nevertheless, a quick Google Scholar search will reveal that decision tree research is not dead. Researchers still continue to propose new or improved algorithms and analyses.

Theory is one thing, but the algorithms need to be implemented as computer programs to actually be useful. Scikit-learn [PVG⁺11] is a very popular machine learning library written in Python. As such, it also contains implementations of various decision tree induction algorithms. Before scikit-learn became popular, a Java-based library called Weka [EHWP16] (or “Waikato Environment for Knowledge Analysis” in full) was often used instead. The implementations of decision tree algorithms in Weka are to this day still in many respects superior to those in scikit-learn. Other tools that implement similar algorithms exist (e.g., Apache Spark [ZXW⁺16]), but those are beyond the scope of this text.

1.2 Goal

The goal of this thesis is to alleviate the discrepancies between scikit-learn and Weka concerning decision tree induction. Mind that decision tree induction tools can never be truly “complete” as stated in the title because the field is immensely broad and still continues to grow. Nevertheless, an effort can be made to improve feature parity between these two popular tools.

One such discrepancy was found when comparing the performance of decision tree induction algorithms in Weka and scikit-learn on an activity dataset [KWM11]. The difference between classification accuracies in this case was considerable at about 25% in favour of Weka.

1.3 Motivation

Some would perhaps question the relevance of such “outdated” techniques anno 2018. This feeling is misguided. The advantages of decision tree induction algorithms are still hard to compete with, even for more modern algorithms [PVG⁺11, Mur98, KZP07]:

1. Comprehensible: makes intuitive sense even for the uninitiated.
2. Transparent, as opposed to for example artificial neural networks
3. Easy to visualize (if number of nodes remains small)
4. Non-parametric: makes very few assumptions about data
5. No data normalization required
6. Handles both categorical and numerical data
7. Handles missing data elegantly
8. Handles multiclass, multilabel and multioutput problems natively
9. Short fitting duration
10. Short prediction duration

Of course decision tree induction algorithms are not perfect:

1. Unstable: small modifications in training set can result in a completely different tree
2. Learning optimal trees is an NP-Complete problem [HR76], so heuristics are used to find approximations
3. Prone to overfitting if not actively countered by adding early stopping criteria or an extra pruning step
4. Prone to bias when one class appears more much frequently in the training set than others.

Some of these drawbacks can be overcome by using an ensemble of decision trees, but that in turn negatively impacts some of the advantages.

1.4 Thesis structure

The structure of the remainder of this text is as follows. First, an overview of the literature study concerning decision tree induction will be presented and the scope of the thesis will be determined in [chapter 2](#). Next, in [chapter 3](#) the decision tree implementations in Weka and scikit-learn are compared to their underlying algorithms and to each other. This results in a list of capabilities. Based on these differences in capabilities, we formulate some hypotheses that can explain the differences in performance. One obvious problem — a missing capability of scikit-learn that does not require additional verification — is attended to in [chapter 4](#). Chapter [5](#) details experiments used to verify whether other discrepancies in capabilities also have a performance impact. It is followed by [chapter 6](#) which presents and discusses the results of these experiments. Additionally, the work done in [chapter 4](#) is validated. We conclude with [chapter 7](#).

Chapter 2

Literature review

The relevant literature for this thesis mostly consists of papers concerning decision tree induction. These go back many decades, but fortunately there are some review and survey papers that provide a convenient overview [Mur98, RM05, KZP07]. On top of the academic literature, the source code and accompanying documentation of scikit-learn and Weka have also been a rich source of information.

2.1 Prerequisites

The reader ought to be familiar with basic machine learning concepts such as supervised learning, classification, regression, overfitting, model validation and ensemble learning. Furthermore, elementary knowledge of decision tree induction is expected. The most important basic concepts will briefly be recapitulated. Topics that are particularly important for the next chapters will be elaborated on.

2.2 Scope

A wide variety of decision tree induction algorithms exists. Here, only the *top down induction of decision trees (TDIDT)* family is considered. It is the most common approach and it is particularly relevant to the software tools under scrutiny.

Unsupervised and semi-supervised algorithms are out of scope. Furthermore, only classification trees are considered. With little effort, most TDIDT classification algorithms can be converted to regression algorithms. Yet, these are far less popular and better alternatives such as XGBoost [CG16] exist.

Ensemble methods are also out of scope. Recent decision tree algorithms rarely work with a single (complex) tree, but rather with an ensemble of relatively simple trees. Random forests [Bre01] is a very popular example of bootstrap aggregating or *bagging*. Regardless, the scope of this thesis concerns the fundamentals of decision trees, and not their derivatives. On the other hand, implementation improvements

suggested in this thesis can potentially benefit related ensemble methods. Other derivatives that are out of scope include logic trees and model trees.

The algorithms in scope are all offline learning methods. This is also known as batch learning. It implies that all computation is done in one place and that all data has to fit in memory. Online learning has gained relevance in the big data era since it does not impose such restrictions. On the other hand, the technique is only used if the context requires it because the performance is not up to par with offline learning.

Finally, only univariate tests are in scope. The test performed in each internal node must only evaluate one attribute of the observation. For categorical attributes, this typically implies checking whether or not the input is equal to a fixed category. For numerical (and thus ordered) attributes, the input value is compared against a fixed threshold using the less than or equal operator. Consequently, the input space is partitioned recursively using axis-aligned hyperplanes. This scope limitation precludes well-known but unpopular extensions such as oblique trees.

2.3 Terminology

Throughout the relevant literature, there is a lack of ubiquitous vocabulary shared by all researchers. Decision trees are used in various scientific fields, each with its own jargon. Specifically, there is a big divide between researchers that approach the problem from a machine learning perspective compared to those who come from a statistics background. To avoid confusion, we avoid synonyms and always use the same term for the same concept. In the following paragraphs, some basic terms are reviewed.

A *decision tree* consists of (*internal*) *nodes* which are connected to other nodes via a one-to-many *parent-child* relation on one hand, and *leaves* which have no children on the other hand. The *root node* is the only node without parent. In a *binary tree*, all internal nodes have two children.

Induction algorithms typically receive a *training set* as input data to construct a decision tree while a *test set* is used afterwards for model validation. These sets are tables of data where each row represents an *observation*. All observation are fully described by a common set of *attributes*. Some attributes are *categorical*, others may be *numerical*.¹ In a supervised learning context, one or more *target attributes* are present. For classification tasks, target attributes are considered categorical even if they contain numeric data. The distinct values of a target attribute are called the *class labels* of that attribute. A classification task is called *binary classification* if there is exactly one target attribute which has exactly two associated class labels and each observation is associated with exactly one of the two class labels. If multiple

¹The latter is sometimes also referred to as *ordered* because this is the underlying property of numbers the algorithm will exploit at some point. Strictly speaking categories can also have an implicit order. In that case, a good practice is to make this explicit by numerically encoding them beforehand.

class labels can be associated with an observation, the task is called *multilabel classification*. If instead there are more than two class labels, it is referred to as *multiclass classification*. *Multioutput classification* on the other hand occurs when multiple target attributes, each with their own distinct set of class labels, have to be derived from the same set of (non-target) attributes. This can be accomplished trivially by creating multiple trees, each handling one class. However, combining them in one tree might offer performance benefits. Decision trees are one of the few machine learning algorithms that can handle all these modes of operation natively.

During *fitting*, first one root node is created and all observations in the training set are stored in this node. When a node is *split* using some *test function*, this function partitions the observations in subsets and then creates a child node for each subset. This process is repeated recursively until some stopping criterion is reached. The *purity* of a node is defined as the percentage of observations in that node that belong to the majority class. A *pure node* is a node with 100% purity. When the tree is complete, it can *predict* the output based on input from, for example, the test set.

2.4 A generic TDIDT algorithm

A typical TDIDT algorithm for classification consists of two phases: a grow phase and an optional prune phase. The grow phase requires three subroutines with fixed signatures: a test generation subroutine, a node splitting subroutine and a stopping subroutine. Historically, researchers presented their TDIDT algorithms with fixed subroutines. Because of the shared interface it is now common to choose these subroutines *à la carte*. One could try to evaluate the performance of each subroutine separately, but selecting the best subroutine for each slot separately does not guarantee a global optimum. Holistic tests must be performed to ensure the best configuration is chosen. Also note that the efficacy of each combination seems to depend on the domain in which it is applied [Min89].

2.4.1 Univariate test generation

Based on the set of observations in a node, tests can be devised that partition this set in a number of subsets. The goal of this step is to generate a finite number of tests $\tau_i \in \mathcal{T}$ based on one associated attribute. Recall the tests based on multiple attributes exist but are out of scope. In the next step, one specific test is chosen from this set of possible tests.

Generating tests for categorical attributes is trivial. For binary trees, each test compares its input (i.e., the value of the associated attribute of an observation) against one fixed value selected from all known values of the associated attribute. If the test is positive, the observation belongs to the first subset, else to the second. This results in as many tests as there are possible values for the associated attribute. For non-binary trees, one test suffices that maps each distinct attribute value to a specific subset. Figures 2.1 and 2.2 illustrate the difference. Non-binary trees are

typically shallower and thus faster to go through from top to bottom during the inference phase. On the other hand, it is an all-or-nothing approach. If an attributes contains 100 different categories where only a few of them are very relevant, is it worth it to add 100 child nodes to the current node? Binary trees have more flexibility in this regard, but they pay for it in depth.

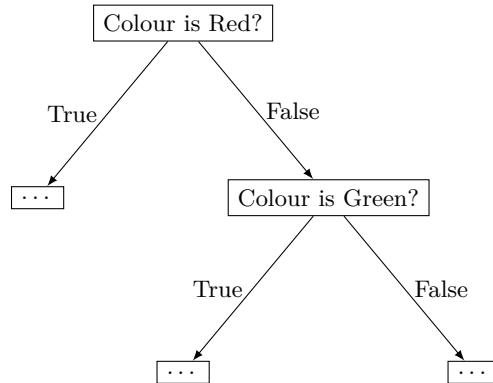


FIGURE 2.1: Example of nodes that perform binary tests associated with the attribute Colour. Two levels are needed to reach the same conclusion as in figure 2.2. Of course in practice not every value is equally relevant, so it is difficult to make statements about which approach is better without knowing the context.

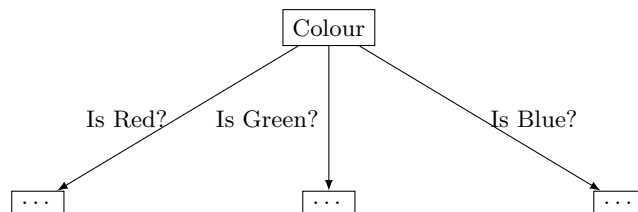


FIGURE 2.2: Example of a node that performs a test associated with the attribute Colour. Since this target attributes contains three distinct class labels, the observations are partitioned over three child nodes.

In the case of numerical attributes, thresholds are introduced to partition the observations based on the implied ordering. That way, an infinite number of tests can be generated, which is of course undesirable. However, at least for the training set, not all tests will result in a different partitioning. A clever choice of thresholds should bring the number of subsets back to a manageable level. Figure 2.4.1 shows an example.

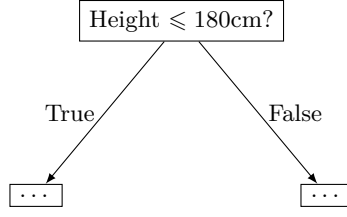


FIGURE 2.3: Example of a node that performs a test associated with the numeric attribute Height. In this case a threshold of 180cm was chosen, but of course any threshold is allowed as long as it divides the set of all height values in two non-empty subsets.

2.4.2 Splitting

Classic TDIDT algorithms work by recursively splitting nodes based on some optimal test $\tau \in \mathcal{T}$, the set of all possible tests. A heuristic called the splitting criterion is required to determine this τ . A few such criteria have stood the test of time.

Purity

The perfect test τ^* creates a partition $\mathcal{S}_{\tau^*} = \{S_1, \dots, S_k\}$ wherein each subset is pure, so optimizing for weighted average partition purity is a sensible first criterion.

$$p(\mathcal{S}_{\tau}) = \sum_i \frac{|S_i|}{|S|} p(S_i) \quad (2.1)$$

Here, $S = S_1 \cup \dots \cup S_k$ and $p(S)$ is the set purity as described above.

Entropy and information gain

In practice purity does not appear to work very well. That is why researchers came up with an alternative based on Shannon's information theory [Sha48]. Quinlan used such metrics in many of his prominent algorithms such as ID3 and C4.5 [Qui86, Qui93], but it was already invented earlier for the Concept Learning System (CLS) [HMS66]. Define entropy (or missing information) of a variable V with possible values v_i and associated probabilities p_i as follows:

$$s(V) = - \sum_i p_i \log_2(p_i) \quad (2.2)$$

The same concept can be applied to the class variable. Define the class entropy $s_C(S)$:

$$s_C(S) = - \sum_c p(c) \log_2(p(c)) \quad (2.3)$$

where $p(c)$ is the probability that a random observation in S belongs to class c . This value can be defined for any node, independent of any specific partition.

For a given test τ , a similar definition can be given for each subset S_i of the induced partition on S :

$$s_C(S_i) = - \sum_c p_i(c) \log_2(p_i(c)) \quad (2.4)$$

For the entropy of the whole partition \mathcal{S}_τ , again use the weighted average entropy of its subsets:

$$s_C(\mathcal{S}_\tau) = \sum_i \frac{|S_i|}{|S|} s_C(S_i) \quad (2.5)$$

Finally, calculate the information gain $h_{IG}(\tau, S)$ of the split that resulted from test τ :

$$h_{IG}(\tau, S) = s_C(S) - s_C(\mathcal{S}_\tau) \quad (2.6)$$

where \mathcal{S}_τ is the partition resulting from test τ .

Gain ratio

The information gain criterion is biased towards tests with many possible outcomes. This could be a problem in non-binary trees. The gain ratio alleviates this problem. First define split information $SI(\tau, S)$ — the maximum possible information gain — as follows:

$$SI(\tau, S) = - \sum_i \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|} \quad (2.7)$$

Finally, define the gain ratio:

$$h_{GR}(\tau, S) = \frac{h_{IG}(\tau, S)}{SI(\tau, S)} \quad (2.8)$$

In binary trees, this heuristic typically causes a less balanced tree compared to the information gain criterion [Qui93].

Gini

Distance metrics such as the Gini impurity index can be used instead of heuristics based on information theory [BFSO84]. The definitions follow the same pattern as those of the information gain:

$$g(S) = \sum_c p(c)(1 - p(c)) \quad (2.9)$$

$$g(S_i) = \sum_c p_i(c)(1 - p_i(c)) \quad (2.10)$$

$$g(\mathcal{S}_\tau) = \sum_i \frac{|S_i|}{|S|} g(S_i) \quad (2.11)$$

$$h_G(\tau, S) = g(S) - g(\mathcal{S}_\tau) \quad (2.12)$$

2.4.3 Early stopping

Breiman argues in his CART book [BFSO84] that choosing good early stopping criteria is far more important than choosing good splitting criteria. If early stopping was not applied or no pruning (see below) was performed afterwards, trees would grow excessively large on real world data sets. This is a classic case of overfitting. It negatively impacts many factors that make decision trees attractive in the first place, such as their comprehensibility and their quick fitting and predicting. It is also detrimental to the performance of the model on unseen data since the model fails to generalize properly.

There are simple and more complex stopping criteria. Simple ones are based on statistics such as these:

1. tree depth
2. number of leaves
3. number of observations in a node
4. purity of a node

More complex stopping criteria can be based on the Minimum Description Length (MDL) of a tree [Ris78] or on statistical techniques such as a χ^2 -test. Quinlan proposed to use the latter in his ID3 algorithm but decided not to include it in the successor C4.5 [Qui86, Qui93] due to its unreliability.

2.4.4 Pruning

A better alternative to early stopping criteria is to let the tree grow freely, and to prune it afterwards in a bottom-up fashion. Typically, the current error estimate of the subtree rooted at the given node is compared to what the estimated error would be if this node was converted to a leaf by pruning away its descendants. If it would perform better as a leaf, the descendants are effectively pruned away. Many different pruning algorithms exist. What follows is a non-exhaustive list of common pruning approaches.

Reduced Error Pruning

Reduced Error Pruning (REP) is one of the most straightforward and statistically sound methods of pruning a tree [Qui87, EK01]. Instead of using the whole training set to grow the tree, some randomly chosen observations are withheld in a separate validation set. By using this validation set after the growth phase is completed, an unbiased estimate of the error of each node in the tree can be calculated. Nodes at the bottom of the tree are converted into leaves if the estimated error of the leaf is equal to or less than the estimated error of the subtree rooted at the given node. This process is repeated recursively until the smallest possible tree is obtained with the minimum estimated error based on the validation set.

The disadvantage of this method is that less data is available for growing the tree, potentially negatively impacting this process. This is not a concern if training data is available in abundance.

Error Based Pruning

Error Based Pruning (EBP) is a technique used in C4.5 [Qui93]. It does not require a separate validation set, so the full training set can be used to grow the tree. The downside of this is that this method is less statistically sound because it is impossible to calculate an unbiased estimate of the error. Instead, an upper bound is calculated based on the training error and that upper bound is used instead of the original error in comparisons. Generally speaking, if a node is associated with fewer observations then there is less certainty about the error and the upper bound will be further away from the original value. The calculation also depends for a large part on the required level of confidence, which the user must provide.

Cost Complexity Pruning

Cost Complexity Pruning, used in the CART algorithm [BFSO84], takes another approach akin to regularization in classic optimization problems. First, it generates a series of pruned trees based on the original. Then it considers both the total training error and a cost factor proportional to the size of each tree to make a first selection. If the training error increases due to the pruning, but it is compensated for by a much smaller tree, the operation as a whole can still be considered positive depending on a trade-off factor. The final tree is chosen from this first selection using

a separate validation set. As such, the same drawbacks apply here as for Reduced Error Pruning.

Others

Many other pruning algorithms exist [Min89, BA97, Elo99, EMSK97]. The reader can find some inspiration in the following list:

1. Minimum error pruning [NB86]
2. Pessimistic pruning [Man97, Qui87, Qui93]
3. MDL-based pruning [MRA⁺95, QR89]
4. Critical Value Pruning [Min87]
5. Pruning using back propagation [KC01]

Alternative: Rule-based Pruning

An outlier in this list is Rule-based Pruning. Decision trees can be converted to a series of if-then statements where the condition is a conjunctive clause. These statements can be further simplified to if-then-else statements and then optimized, which can be seen as an alternative form of pruning. The resulting model is no longer a tree, so it is considered out of scope for this thesis, but it can still approximate the underlying concept that the tree used to represent.

2.5 Conclusion

TDIDT algorithms incorporate different subroutines, for each of which a number of alternatives are available. This makes them a very flexible tool with uses in a variety of settings. Popular algorithms such as C4.5 and CART are opinionated in the sense that they each propose a small number of specific configurations of components. Fortunately, that does not stop algorithm implementers from offering more choice to their users, as shown in the next chapter. Note also that there is no single precise definition of ID3, C4.5 or CART. New insights were acquired over time and added to the solution, but the algorithm name rarely changed.

Chapter 3

Existing implementations

Chapter 2 discussed the theoretical basics of decision trees. In this chapter, we consider applications of this theory in the form of two popular software libraries: Weka [EHWP16] and scikit-learn [PVG⁺11]. The focus is exclusively on Weka’s J48 algorithm and scikit-learn’s `DecisionTreeClassifier`. Scikit-learn only offers one other decision tree induction algorithm that performs regression instead, but that is out of scope. Weka on the other hand contains a variety of other tree builders, but most are derivatives or ensemble forms of the original J48 algorithm. Model trees and an online learner are also included, both beyond the scope of this thesis.

3.1 Capabilities

The most important difference between the two implementations regarding decision trees is in the base algorithm they started from. The J48 algorithm in Weka is based on Quinlan’s C4.5 [Qui93], while the `DecisionTreeClassifier` in scikit-learn used CART by Breiman [BFSO84] as the foundation. This has a profound impact on the capabilities of both implementations. These capabilities are divided in categories and discussed one by one in the following subsections.

3.1.1 Structural capabilities

CART only supports binary trees, and the same applies for scikit-learn’s implementation. It is an option in J48, but the default settings generate non-binary trees. Binary trees are typically deeper than their non-binary counterparts, making the inference phase more computationally expensive. On the other hand, Elomaa et al. claim that the use of binary discretization with C4.5 needs about the half fitting time of using C4.5 multisplitting [ER99]. Note that this only impacts splits on categorical attributes; splits of numerical attributes are always binary.

3.1.2 Input capabilities

Categorical and numerical attributes

ID3 was not capable of dealing with numerical attributes, but C4.5 and CART can handle both types. J48 can also handle both just like its theoretical counterpart C4.5. Scikit-learn’s implementation however did not inherit the categorical input capabilities of CART. This is understandable from a software engineering perspective: scikit-learn is built on top of NumPy [Oli06] which itself is a numerical, scientific computing library. The user can work around this issue by pre-processing the categorical data, using for example a `LabelEncoder`¹ to map each category onto a unique ID or a `OneHotEncoder`² that introduces as many boolean dummy variables as there are categories where only one is active at any given category.

Consider an example where an attribute Colour contains three categories: Red, Green and Blue. By default C4.5 and Weka would generate a single test out of this attribute. Red would be mapped to the first child, Green to the second and Blue to the third (cf. Figure 2.2). CART on the other hand only generates binary trees, so it would generate three tests instead: “Colour is Red?”, “Colour is Green?” and “Colour is Blue?”. Each test has a boolean output which decides whether to jump to the left or the right child. Only two of the three tests need to be used to determine the colour: if the colour is neither Red nor Green, it is considered Blue (cf. Figure 2.1). More generally, for N distinct categories in an attribute, $N - 1$ tests have to be performed to cover all states. Note that not all states necessarily have to be checked in a real world scenario.

In scikit-learn, a one-hot encoder is needed before the data can be processed. It would replace the Colour attribute with three numerical dummy attributes: ColourRed, ColourGreen and ColourBlue. In this case, Red is represented as $[1, 0, 0]$, Green as $[0, 1, 0]$ and Blue as $[0, 0, 1]$. Because these are numerical attributes, the test generation function will try to find thresholds to partition the space. In this case, only one threshold somewhere between 0 and 1 (e.g., 0 itself) is needed. As such, each dummy variable will cause the generation of one test: $(x \leq 0)$. Although the implementation is slightly different, there is a trivial isomorphism between these three tests and the three CART tests. Consequently, the semantics are preserved even though they are slightly obscured. Compare Figure 2.1 from the previous chapter with Figure 3.1 to confirm this. On a hardware level, there is no intrinsic difference in performance between the $=$ and the \leq operator [Fog16].

Alternatively, the label encoder would replace this categorical attribute with a single numerical attribute with possible values 0, 1 and 2. This implies that there is an order among the colours, which is not supposed to be the case. Tests such as $(x \leq 1.5)$ could be generated. This is equivalent to “Colour is Red or Green?”, which

¹<http://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.LabelEncoder.html>

²<http://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.OneHotEncoder.html>

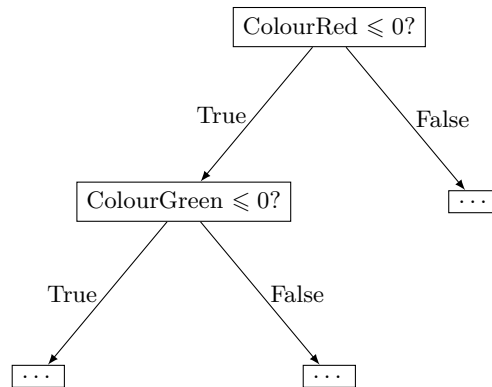


FIGURE 3.1: Reprise of the colour example where the data has been pre-processed using a one-hot encoding. Compared to [Figure 2.1](#) the structure is mirrored and the node names are less clear, but the essence remains the same.

is more expressive than what was possible so far. That looks positive at first, but not all logical combinations can be represented this way. In this example, no threshold can be found that implies “Colour is Red or Blue”. Even worse, what can be represented depends entirely on the implied yet non-existent order of the original categorical variable. In short, this technique does not adhere to the original semantics and should not be used.

In conclusion, the suggested workaround with the one-hot encoder is acceptable but that does not change the fact that it negatively impacts two of the decision tree advantages we listed earlier in [chapter 1](#): no data preparation required and excellent comprehensibility.

Missing values

Another important input capability is dealing with missing values. As stated before, decision trees have fairly straightforward ways of dealing with this problem. While the papers on ID3 ignored this issue, both C4.5 and CART have methods for dealing with it. As before, Weka inherited the capability from C4.5, but scikit-learn did not inherit the same from CART. This is another case where extra pre-processing is needed to make decision tree induction work with scikit-learn. During this process, valuable observations might be thrown out entirely.

3.1.3 Output capabilities

Classification and regression

For this thesis regression is out of scope, but it is still an important point. CART has both classification and regression in its name, so it comes as no surprise that both are supported by this algorithm. The scikit-learn implementation is split in

a `DecisionTreeClassifier` and a `DecisionTreeRegressor`. C4.5 only supports classification out of the box. Consequently, regression is not a capability of Weka's J48 implementation. However, Weka contains other decision tree algorithms, some of which do accommodate regression (e.g., `REPTree`).

Multiclass, multilabel, multioutput

On this front, scikit-learn is the superior implementation because all variants of multiclass, multilabel and multioutput are supported. Weka's J48 implementation on the other hand only supports multiclass. CART and C4.5 have the same limitations as Weka's. A simple workaround for multioutput is to fit multiple trees instead, although this might impact various performance aspects.

3.1.4 Splitting criteria

Another difference lies in the choice of splitting criteria. C4.5 and J48 support the typical criteria based on information theory such as information gain and gain ratio. CART and scikit-learn on the other hand both support information gain and the Gini impurity index but not gain ratio since it is not beneficial for binary trees. The purity criterion is not implemented in any of the algorithms due to its poor performance.

3.1.5 Early stopping and pruning

As mentioned at the end of [chapter 2](#), both C4.5 and CART seemingly have their favourite pruning algorithms. The former explicitly supports Error Based Pruning and Rule-based Pruning, while the latter favours Cost Complexity Pruning. J48, like C4.5, supports Error Based Pruning by default. Quinlan's other proposal, Reduced Error Pruning, is also an option. Rule-based pruning is not supported, but Weka contains other rule based algorithms as an alternative. Of course the option not to prune at all is also available. In that case users can fall back on early stopping criteria. All algorithms and implementations include some simple stopping criteria. ID3 also introduced more complex stopping criteria, but that effort was abandoned in favour of pruning in C4.5. Breiman came to the same conclusion for CART. Surprisingly, scikit-learn did not follow Breiman's suggestion of implementing Cost Complexity Pruning. Instead, it opted for the inferior practice of using simplistic early stopping criteria such as a maximum tree depth or a minimum number of observations per node required to split it. This is by far the most significant difference between the two implementations.

3.1.6 Mode of operation

Machine learning algorithms typically operate in either one of two modes: online or offline (batch) mode. None of the algorithms and implementations under scrutiny offer online learning. Weka does include an implementation of an online tree builder (i.e., a `HoeffdingTree`), but not as part of J48.

3.1.7 Miscellaneous

Many more capabilities exist, but these are beyond the scope of this thesis. The interested reader can dive deeper into any of the following topics:

- Generating rulesets starting for a tree representation
- Assigning different weights to different classes
- Assigning different weights to different observations
- Taking asymmetric misclassification costs into account

3.1.8 Summary

Figure 3.2 shows a summary of the previous paragraphs. During the comparison of capabilities of the C4.5 and CART algorithms and their respective software counterparts Weka J48 and scikit-learn's `DecisionTreeClassifier`, three important discrepancies at the expense of the latter came up.

- It lacks pruning algorithms, making it rely on inferior simple stopping criteria instead.
- It cannot handle categorical attributes natively. This specifically impacts its capability of making categorical multisplits.
- It fails when given a training set with missing data.

Other discrepancies came up, but these were either minor or could be rationalized as explained in the previous paragraphs.

3.2 Next steps

3.2.1 Lack of pruning

The lack of pruning is a particularly noteworthy discrepancy. The current workaround involves carefully defining early stopping criteria. However, it is clearly established in the literature that this workaround is inferior to the usage of proper pruning algorithms [BFSO84, Qui93]. Fixing this is priority number one. Fixing in this sense means implementing pruning algorithms for the scikit-learn decision tree induction algorithms. The next chapter immediately discusses this in more detail. In later chapters, we propose an experimental setup to evaluate the performance of the new implementation and present the results of that evaluation. Our hypothesis: pruned trees take longer to train but are smaller. That makes them less prone to overfitting so the accuracy should be the same or better and the prediction time should decrease.

3. EXISTING IMPLEMENTATIONS

▼ Capability	▼ ID3	▼ C4.5	▼ J48	▼ CART	▼ DTC	▼
Structure						
Binary numeric splits	N	Y	Y	Y	Y	
Binary categorical splits	N	Y	Y	Y	Y	
Categorical multisplits	Y	Y	Y	N	N	
Input						
Numerical attributes	N	Y	Y	Y	Y	
Categorical attributes	Y	Y	Y	Y	N	
Missing values	N	Y	Y	Y	N	
Output						
Binary classification	Y	Y	Y	Y	Y	
Multiclass classification	N	Y	Y	Y	Y	
Multilabel classification	N	N	N	N	Y	
Multioutput classification	N	N	N	N	Y	
Regression	N	N	N	Y	Y	
Model trees?	N	N	N	N	N	
Splitting						
Purity	N	N	N	N	N	
Information gain	Y	Y	Y	Y	Y	
Gain ratio	N	Y	Y	N	N	
Gini impurity index	N	N	N	Y	Y	
Early stopping						
Simple	Y	Y	Y	Y	Y	
Complex (χ^2 , MDL, ...)	Y	N	N	N	N	
Pruning						
None	Y	Y	Y	Y	Y	
Reduced Error Pruning (REP)	N	N	Y	N	N	
Error Based Pruning (EBP)	N	Y	Y	N	N	
Cost Complexity Pruning (CCP)	N	N	N	Y	N	
Rule-based pruning	N	Y	N	N	N	
Mode						
Streaming	N	N	N	N	N	
Batch	Y	Y	Y	Y	Y	
Misc						
Generate rulesets	N	Y	N	N	N	
Class weights	N	N	Y	Y	Y	

FIGURE 3.2: Summary of various capabilities as seen in the ID3, C4.5 and CART algorithms, but also in their implementations: Weka's J48 and scikit-learn's DecisionTreeClassifier (DTC).

3.2.2 Categorical attribute handling

Next to the lack of pruning, scikit-learn's algorithms also have a problem with categorical input. Because scikit-learn is built on top of NumPy, it uses numerical matrices to pass data between all its components. As such, categorical data simply cannot serve as input to most of its algorithms. It is part of the core philosophy of scikit-learn. Fixing this head-on would require a complete rewrite of not only the decision tree algorithms, but also all the basic interfaces and their dependencies, which is far beyond the scope of this thesis. On top of that, the use of NumPy is essential for the performance of scikit-learn. Python in itself is a slow, interpreted scripting language that is not suited for high performance computing tasks. Its integration options with C code (and the intermediate Cython code) make it a viable option, but at a considerable development cost. NumPy abstracts this complexity away behind a pure, easy to use Python interface while maintaining the speed boost.

A less significant but related discrepancy is that scikit-learn's decision tree induction implementation only generates binary trees. Breiman made this conscious choice for his CART algorithm in 1984, and it remains so today in its derivatives. Fixing this once again fundamentally changes the implementation and therefore entails a complete rewrite of the decision tree package.

In conclusion, the currently accepted workaround of pre-processing the categorical attributes in the input data using a one-hot encoder seems to be our best bet. In the upcoming chapters, we design an experiment that tests whether this workaround has any performance implications. Given the isomorphism between binary tests based on categorical attributes and binary tests based on numerical values that are the result of a one-hot encoding, we hypothesize that the performance impact is limited.

3.2.3 Missing value handling

The third and lowest priority problem is the failure of scikit-learn to process datasets containing missing values. Three approaches are possible. The first workaround discards any observation with missing values. Depending on the context, this data loss is somewhere between irrelevant and unacceptable. The second workaround attempts to guess the values for the gaps in the data based on the rest of the data. The user must judge whether this approach is feasible in his context, but it certainly is no universal solution. For that, we look at the third option. Unlike the previous workarounds, this is not a pre-processing step. Instead, the classifier itself must be modified to deal with missing data directly. Decision tree induction algorithms typically have that capability. The C4.5 and CART texts both outline ways to deal with this problem. One option is to split an observation with a missing value into as many fake observations as there are possible values for the missing (categorical) attribute where each distinct value is used once. Some lower internal weight could be assigned to these fake observations to compensate for the data expansion.

3.3 Conclusion

This chapter introduced us to the various capabilities of decision trees. Those capabilities were mapped onto both the CART and C4.5 algorithms, but also the Weka J48 and scikit-learn implementations. This exercise revealed three important discrepancies that make scikit-learn inferior to Weka regarding decision tree induction. Those discrepancies concern pruning, categorical values and missing values. For each of these, next steps have been defined. The next chapter immediately tackles the pruning problem head-on.

Chapter 4

Implementation extensions

In the previous chapter, we found three important features that scikit-learn lacks in comparison to Weka’s J48 implementation. In this chapter, we tackle the high priority pruning issue.

4.1 New pruning functionality

The most straightforward way to solve the pruning problem in scikit-learn is simply to implement an extension that provides this feature. This extension is simply a python package that can be downloaded and installed by anyone who already installed scikit-learn beforehand. We developed a new classifier called the `PruneableDecisionTreeClassifier` as the main component of such an extension.

Instead of developing an extension, another option was to fork the complete current scikit-learn codebase and modified that in-place instead. The advantage of this approach is that pull requests can be made on Github to update the original codebase. However, in that case the code must adhere to a very strict contribution policy with requirements beyond the scope of this thesis. For example, scikit-learn must work perfectly with both Python 2.x and Python 3.x distributions. In this thesis we follow the way forward and only support Python 3.x. Additionally, keeping the code up to date involves regular merges with new scikit-learn code. This is not a problem with extensions built against a specific scikit-learn version (i.e., 0.19.1).

4.1.1 Functional requirements

The new class must enable tree pruning, but it must also function in a way highly similar to the original decision tree classifier in scikit-learn. This means it must:

- have at least the same constructor arguments such as `criterion` and `splitter`
- support various pruning strategies

- have additional constructor arguments to specify the pruning method and related settings
- adhere to the basic classifier interface containing methods such as `fit`, `predict`, `predict_proba` and `score`
- implement typical tree methods such as `apply` and `decision_path`
- offer typical tree attributes such as `classes_`, `feature_importances_` and `tree_`

These requirements facilitate migrations to the new pruning-enabled decision tree classifier. Additionally, the new classifier remains compatible with existing scikit-learn tooling such as pipelines and grid searches.

From a software engineering perspective, meeting these requirements can be accomplished in two ways. Either make `PruneableDecisionTreeClassifier` inherit from `DecisionTreeClassifier` or implement the former using a decorator pattern. The result is the same, but the second option would require more boilerplate code. As such, the first option (inheritance) was chosen. The class diagram is shown in figure 4.1.

Adding pruning to an existing algorithm exclusively involves adding another step to the fitting process. In scikit-learn terms, this only requires an override of the `fit` method, while the other methods keep their existing implementation.

4.1.2 Pruning algorithms

Various pruning strategies exist, as discussed in chapter 2. The goal is feature parity with J48, so both Reduced Error Pruning and Error Based Pruning are implemented. Additionally, the algorithm offers the option to disable pruning by setting the constructor argument `prune=None`. This way, users have no need to import the original classifier anymore.

Implementation-wise, the pruning logic has been isolated from the other tree logic in a separate class hierarchy. This is illustrated in figure 4.2. A base class `Pruner` is provided that contains all the tree operations that pruning algorithms will need such as `is_leaf`, `to_leaf` and `leaf_prediction`. Two subclasses are also provided, the `ReducerErrorPruner` and the `ErrorBasedPruner`. Each one exposes a `prune` method. This way, supporting a new pruning strategy is as simple as adding a new class to this hierarchy, instantiating it from the classifier and calling the `prune` method after the growth phase.

Reduced Error Pruning

Reduced Error Pruning is enabled by setting the constructor argument `prune='rep'`. It requires a separate validation set. There are two ways to deal with this problem.

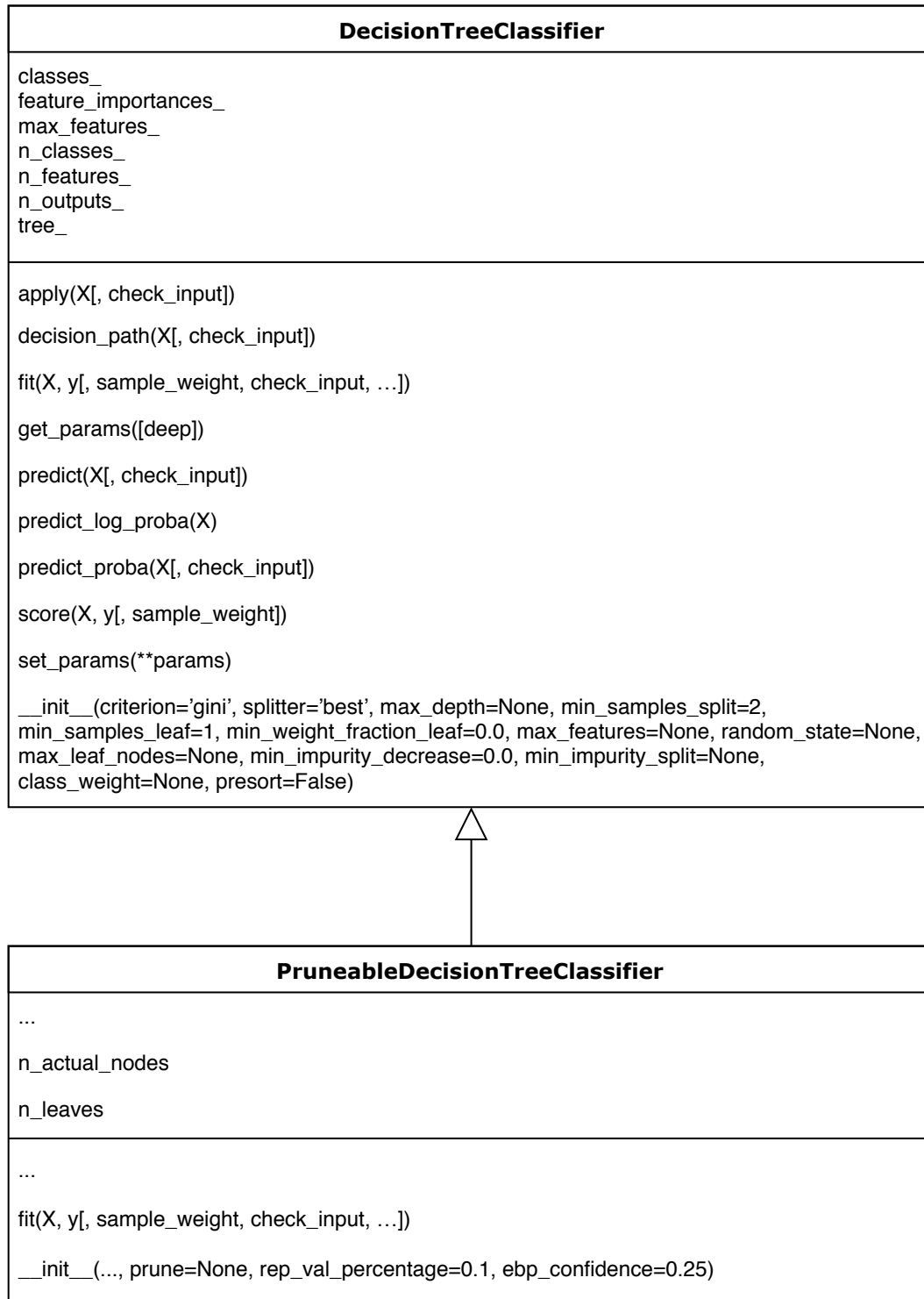


FIGURE 4.1: Class diagram of decision tree classifiers.

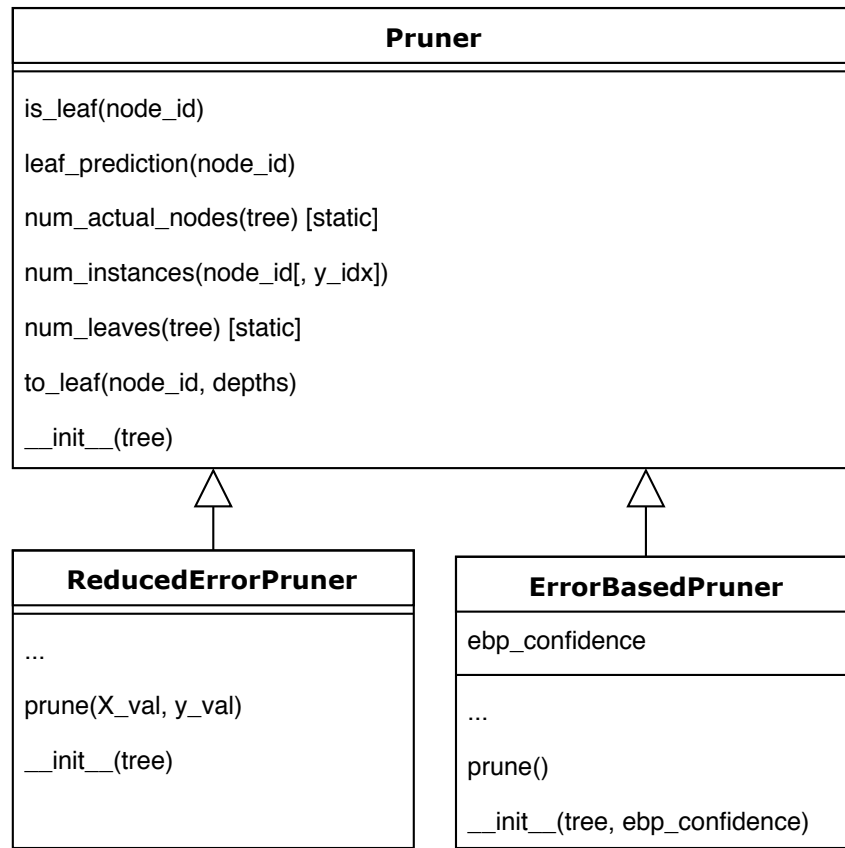


FIGURE 4.2: Class diagram of pruners.

One, introduce two additional parameters to the `fit` method where the user can supply his own validation data (`X_val`) and labels (`y_val`). This maximizes the flexibility. However, it also places the burden of separating the original training set on the user. Worse, it breaks the classifier interface. Option two instead takes care of the separation internally using a stratified `test_train_split`. The user only has to provide the training set as usual, and also indicate which percentage of that set must be reserved for pruning via the `rep_val_percentage` argument. Because interface consistency is part of the functional requirements, option two was chosen. Weka's J48 contains a similar solution when this pruning strategy is used.

Once the two sets are separated, fitting can begin with the reduced training set using the existing implementation (i.e., the `fit` method of the base class). Afterwards, the `prune` method of a `ReducedErrorPruner` instance is called to simplify the tree.

Error Based Pruning

Error Based Pruning works similarly by setting the constructor argument `prune='ebp'`, but it does not require a validation set. Consequently, the training set provided by the

user can be passed directly to the original `fit` method of the `DecisionTreeClassifier` base instance. Next, an instance of `ErrorBasedPruner` is created and its `prune` method is called with a confidence value that the user provided in the constructor of the new `PruneableDecisionTreeClassifier`. This confidence value is needed to calculate a statistical upper bound to the observed errors.

4.2 Supporting tools

A lack of tree pruning algorithms was not the only problem of scikit-learn when compared to Weka's J48. It also has problems with categorical attributes and missing values. We will look at these problems in more detail later on. For now, some straightforward workarounds are implemented to make it possible to properly benchmark the new pruning solution. To that end, a `CsvImporter` class has been added to this extension that can parse arbitrary datasets. It adheres to the scikit-learn transformer interface, offering `fit`, `transform` and `fit_transform` methods. The input is a file path to a Comma Separated Value (CSV) file, and the output is an `X`, `y` pair ready for consumption by tree induction algorithms¹.

NumPy and SciPy — the base libraries that scikit-learn is built upon — already offer similar methods, but they are not suitable for our use cases. For example, SciPy offers a method to parse Weka's Attribute-Relation File Format (ARFF). This format contains valuable information, such as whether a variable is numerical or categorical (including the domain). This information is returned by the method as a structured array that keeps data type information per column instead of one data type for all values. Unfortunately these special arrays are not compatible with most machine learning algorithms. Numpy on the other hand offers a method to read plain text files such as CSV files into regular arrays. However, this function is known in the SciPy community to be unreliable and it also lacks modern features that libraries such as Pandas have available.

Our new solution is based on Pandas. Since it starts from plain CSV files instead of ARFF files, it has to infer the data types of each column based on the contents. Next, it deals with missing values rather abruptly by deleting all incomplete observations. A warning is triggered if a given threshold of data loss is surpassed. Scikit-learn offers more advanced methods of dealing with missing values independent of the machine learning algorithm used in the `impute` submodule [ld18]. The simple imputer only looks at other values in the same column to guess the value of a missing entry, which is also a very blunt method. The advanced version on the other hand makes various assumptions about the distribution of the data that cannot be guaranteed. These methods can be interesting for specific datasets, but our solution must work across a wide range of datasets. Consequently, the straightforward dropping strategy remains in use until a decision tree specific algorithm is implemented to alleviate this problem.

¹Strictly speaking, the `fit_transform` method can only return the data `X` but not the labels `y` due to interface restrictions. Instead, `fit_transform_both` was introduced as a workaround.

In the next phase, the class column — as indicated by the user — is separated from the rest of the data and the labels are encoded to integers in the range $[0, n_classes - 1]$. The original values are also preserved to enable an inverse transformation after classification if needed.

Finally, in the rest of the data, values of categorical attributes are encoded in a one-hot fashion to avoid creating an implicit ordering as discussed earlier in [section 3.1.2](#). At this point the dataset contains no more missing or categorical values. It is ready to be used in decision tree induction algorithms.

4.3 Challenges and limitations

Scikit-learn is written in Python, an interpreted language that was not intended for development of high performance software. The language does have a few tricks up its sleeve to alleviate this problem, but it comes at a cost. That cost becomes particularly apparent when trying to extend existing algorithms. The tricks make the code much harder to understand, maintain and adapt. This contrasts strongly with the clean, object oriented development approach of Weka.

For example, whereas the internal data structures Weka uses to represent its decision trees are also trees, scikit-learn stores everything in non-hierarchical arrays with an entry for each node in the tree. Logical operations on trees, such as traversing the structure, splitting or pruning a node are trivial when the underlying data structure matches the represented structure. Using flat arrays, these operations are more challenging. In the `Pruner` base class outlined above, some of these previously missing operations are now also available for any subclass to use. However, the consequences of the flat arrays have not been overcome entirely. The arrays all have a length equal to the number of nodes in the unpruned tree and the index of the node in the array is also its ID number. During the pruning process, when a node is removed, its entries in the various arrays are invalidated but not removed. This means memory usage of the tree is not reduced by the pruning process, even though the actual tree size might have shrunk considerably. Compressing these data structures after pruning is extremely challenging because all node IDs would have to be remapped. Since the nodes exist in both Cython and Python, this becomes very complicated. A better approach might be to reconstruct the final tree from scratch. However, poor software design choices in scikit-learn have made it so that manipulation of the internal data structure is tightly coupled to the actual decision tree induction. Without rewriting a significant part of the code, this option is also off the table. Instead, the arrays remain uncompressed and the memory usage is not improved.

Implementing this functionality in Weka (had it not already been present) would have been a significantly easier. We understand the conundrum that the scikit-learn developers have faced: using plain, clean Python without any *tricks* would have made the software too slow to be useful in a professional setting. However, we speculate that poor design choices such as the ones outlined here are one of the most important

reasons why scikit-learn still lags behind Weka in certain aspects. The codebase has become so complex that even trivial changes take a long time to implement.

4.4 Result

The resulting implementation is publicly hosted at <https://github.com/vhsven/vhsven-sklearn>. Comprehensive documentation and several examples are all available at <https://vhsven.github.io/vhsven-sklearn>. Check out the three examples concerning decision surfaces to get an intuitive feel of what the effect of each pruning method is.

Chapter 5

Methodology

So far, relevant algorithms were covered in [chapter 2](#) and their implementation counterparts were discussed in [chapter 3](#). At that point, three capabilities were identified that made the decision tree induction implementation in scikit-learn inferior to that in Weka. For the highest priority issue (i.e., lack of pruning), the literature offered sufficient argumentation to immediately tackle this issue. The solution was presented in [chapter 4](#). At this point three tasks remain:

- Validate whether the new pruning implementations work as expected
- Verify whether the accepted workaround for categorical attributes has a significant performance impact
- Find a more elegant solution to the missing values problem

At least for the first two tasks, experiments have to be performed. This chapter explains the experimental setup while [chapter 6](#) covers the actual results. In the first section, the pruning validation methodology is covered. The next section discusses the categorical attribute experiment.

5.1 Pruning implementation validation

First, the experimental setup of the pruning validation is discussed. The experiment involves both the new `PruneableDecisionTreeClassifier` and Weka's J48 classifier. The original `DecisionTreeClassifier` is not included because the pruneable version can be configured to behave exactly the same. The following configurations are used, showing the specific scikit-learn and Weka parameters:

- no pruning and no early stopping
- early stopping using $\text{min_samples_leaf} = \frac{0.5}{\text{n_classes}}$ [scikit-learn exclusive]

- REP with `rep_val_percentage = 10%` or `numFolds = 10`
- REP with `rep_val_percentage = 20%` or `numFolds = 5`
- REP with `rep_val_percentage = 50%` or `numFolds = 2`
- EBP with `ebp_confidence = confidenceFactor = 10-5`
- EBP with `ebp_confidence = confidenceFactor = 10-3`
- EBP with `ebp_confidence = confidenceFactor = 10-1`

Comparing only pruned versus unpruned trees does not tell the whole story. The pruned trees must also be compared against a sensible baseline of what scikit-learn users do if no pruning algorithm is available. That is why the early stopping configuration with `min_samples_leaf` is also in this list. Unfortunately, finding a good value of this parameter that works across a wide range of datasets is not easy. The current formula depends on the context to (partly) work around this issue: it takes the number of classes into account. Note that for this specific early stopping configuration a mirror configuration in Weka is absent.

A comparison within scikit-learn of pruned versus unpruned trees with or without early stopping criteria is a good start, but an external benchmark against similarly configured J48 instances is also interesting. Although the base algorithms of the two software libraries differ, the goal is to maintain maximum compatibility in the experimental setup. To that end, some additional rules must be taken into account.

- All algorithms must produce binary trees
- Multilabel and multioutput capabilities must not be used
- Nodes must be split using the information gain criterion
- Early stopping criteria must not be used unless explicitly specified in the configuration
- Advanced J48 features such as tree collapsing, subtree raising or MDL corrections must be disabled

All other parameters are set to their default values.

For both scikit-learn and Weka, a grid search algorithm loops over all the listed configurations, fits them to a training set and finally compares the prediction made for an independent test set against the correct result. Because we start from a complete data set, it has to be split in a training set and a test set. Tenfold stratified cross-validation takes care of this. For statistical robustness, the whole procedure is repeated ten times. Each configuration will thus be tested one hundred times.

Relevant statistics of each test are gathered by the built-in features of both tools. To have reproducible results, the seeds of the random number generators take fixed values before the experiment starts. Furthermore, the whole process is scripted and can be done in one run each for scikit-learn and Weka. A single thread is used to avoid the complexities of multithreading. All tests are performed on a single computer that is otherwise idle (except for system processes). This ensures that the external environment changes as little as possible and that the test results can be compared in a fair way.

This whole procedure assumes one dataset was given as input. To get a more complete picture of the implementation performance, the test is repeated for various datasets. [subsection 5.1.1](#) introduces all datasets used in this thesis. The aforementioned `CsvImporter` ([section 4.2](#)) pre-processes each dataset when a scikit-learn classifier is used. Weka does not require pre-processing, but all missing values are discarded nonetheless to effectuate an unbiased comparison.

Both scikit-learn and Weka offer tools to facilitate such tests and gather the statistics afterwards. Every configuration has a fairly close counterpart in the other tool. The main differences that remain after taking the above setup into account, are the different base algorithms (C4.5 versus CART) and the different runtimes (Java versus Python).

The new implementations of Reduced Error Pruning and Error Based Pruning are judged on three aspects: the accuracy on unseen data and the fitting & prediction time of the algorithm. The size of the tree and the number of leaves are also saved to get more insight into any potential differences.

5.1.1 Datasets

Part of the experimental setup is the choice of datasets. We opt for a variety of real world classification datasets to get a good picture of the pruning implementation's performance. These datasets are primarily taken from www.openml.org [[VvRBT13](#)]. [Table 5.1](#) gives an overview of all datasets.

These specific datasets were selected for their popularity in the machine learning community much more than for their specific contents. Nevertheless, care was taken to have a healthy mix of datasets with some variation in the number of instances, the number of classes, the number of (categorical) attributes and the amount of missing data.

The *activity* dataset is the odd one out. It belongs to a study by Kwapisz et al. [[KWM11](#)]. In this study, the activity of a test subject is predicted based on sensor data from cell phone accelerometers. The possible activities are walking, jogging, going upstairs, going downstairs, sitting and standing. Scikit-learn's regular `DecisionTreeClassifier` performed poorly on this set in the past, although Weka's J48 handled it well.

Name	Description	C	A	N	M	CA
diabetes	Pima Indians diabetes	2	8	768	No	0
ionosphere	Johns Hopkins ionosphere	2	34	351	No	0
iris	Fisher’s iris	3	4	150	No	0
wine	Wine recognition	3	13	178	No	0
wdbc	Breast cancer Wisconsin	2	30	569	No	0
letter	Letter image recognition	26	16	20 000	No	0
houses	House price high/low	2	8	20 640	No	0
heart-statlog	Heart disease	2	13	270	No	0
monks-problems-2	Monks problems	2	6	601	No	6
tic-tac-toe	Tic-tac-toe endgame	2	9	958	No	9
german_credit	German credit risk	2	20	1 000	No	13
vote	1984 US votes	2	16	435	Yes	16
hepatitis	Hepatitis survival	2	19	155	Yes	13
activity	Activity prediction	6	45	5 424	Yes	0

TABLE 5.1: Overview of real world classification datasets. Column C indicates the number of classes, A the number of attributes (excluding the class attribute), N the number of observations, M whether the datasets contains missing values and CA the number of categorical attributes (also excluding class).

5.2 Evaluate the need for categorical attribute support

One of the conclusions in [chapter 3](#) was that scikit-learn does not support categorical attributes out of the box. The proposed workaround consists of encoding all categorical attributes in a dataset using a one-hot scheme during the pre-processing stage. We hypothesized that this step would not negatively affect the performance of the resulting decision tree classifier. Unfortunately this hypothesis cannot be tested directly in scikit-learn, so Weka is used as a surrogate instead. The J48 can be configured in a way that strongly mimics the behaviour of CART, for example by generating only binary trees and disabling advanced features. If the hypothesis holds for J48 in this configuration, it is reasonable to assume that it would also hold for scikit-learn’s decision tree classifier and consequently that there is no need to extend the latter’s implementation.

The methodology is exactly the same as the one described in the previous section, except that missing values are retained and that other classifier configurations are involved. More specifically, there are only two classifiers. One is the regular Weka J48 classifier with Error Based Pruning enabled and a fixed confidence value. The other is a metaclassifier called a `FilteredClassifier` that combines the supervised `NominalToBinary` attribute filter with another J48 instance configured exactly the same as the first. The filter applies a one-hot encoding on all categorical attributes.

The built-in Weka analysis tool determines whether there is a significant difference in accuracy, fitting or prediction time between these two classifiers at confidence level

$\alpha = 0.05$. This tool uses a statistical test called a (corrected) paired t-test. Its null hypothesis H_0 states that the mean difference between two sets of observations equals zero. The reason the test statistic is corrected is to compensate for the bias introduced by the repeated cross-validation procedure [GB18]. Furthermore this test makes some assumptions about the data, for example that it is distributed sufficiently Gaussian-like. A histogram and a QQ-plot indicate that this is a reasonable assumption in our case. More advanced comparison techniques such as nonparametric comparison tests based on rankings are known [Dem06], but are not readily available for us to use.

As a bonus, the same methodology is used to determine whether non-binary trees perform better than binary trees. To accomplish this the metaclassifier configuration is switched out for a copy of the first J48 configuration but with the binary-only option disabled.

5.3 Conclusion

In this chapter we discussed the experimental setup for three experiments:

- Validate whether the new pruning implementations work as expected
- Verify whether the accepted workaround for categorical attributes has a significant performance impact
- Verify whether binary trees perform significantly different from non-binary trees.

The results of these experiments are shown and discussed in [chapter 6](#).

Chapter 6

Results and discussion

In this chapter the results of the experiments detailed in the previous chapter are shown and discussed. The first section deals with pruning-related results. The next section is about categorical attributes and the final section pertains to the binary tree issue.

6.1 Pruning

In [chapter 4](#), pruning extensions were added to scikit-learn’s decision tree implementations. More specifically, Reduced Error Pruning (REP) and Error Based Pruning (EBP) were implemented. In this section we evaluate how these pruning extensions affect the size of the tree (i.e., the number of nodes), the performance in terms of accuracy and the fitting and prediction speed.

6.1.1 Tree size

The mean size of the trees is shown as a heatmap in [Figure 6.1](#) for each dataset and each scikit-learn configuration. The values have been normalized by the mean size per dataset of unpruned trees. As such, the *none* row which indicates that no pruning was performed is filled with 100% values.

One value exceeds 100%. This is possible because the pruned trees are not derivatives of the unpruned trees. They are simply two different configurations that are tested independently of each other. In this specific case, Reduced Error Pruning was used with a validation set size equal to 50% of the training set size. Since 10% of the dataset was already set apart in the cross-validation procedure, only 45% of the original data was left to train the tree. Compared to the other algorithms where larger fractions of training set remain, this can profoundly change the tree structure — including its size — before pruning. This size is unknown, but it must have been larger than the tree size of an unpruned tree trained on 90% of the original data. Even after pruning, the former tree turned out to be larger than the latter.

6. RESULTS AND DISCUSSION

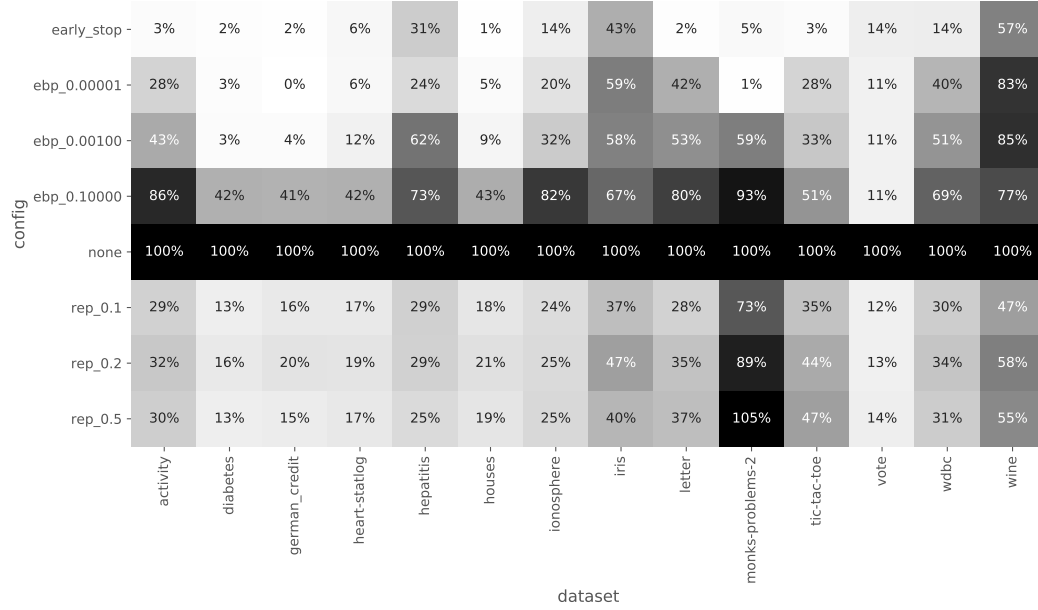


FIGURE 6.1: Mean tree size per scikit-learn configuration and per dataset. The values have been normalized by the mean size per dataset of unpruned trees. Weka results are not included in this heatmap. The colour scale goes from 0% (white) to 100% (black).

The next observation is that the other rows in the heatmap are quite light, meaning that in general the tree sizes have been reduced greatly by the various pruning strategies. For Error Based Pruning, the confidence factor plays a large role in the amount of pruning that is performed. A very low confidence factor of 10^{-5} results in very aggressive pruning while a factor of 10^{-1} results in far less pruning. The validation set size of Reduced Error Pruning on the other hand does not have a clear impact on pruning aggression. A larger validation set only results in more reliable error estimates, but on average these estimates do not move in any particular direction.

Specifically for scikit-learn, a pseudo-pruning configuration is added that uses an early stopping criterion based on the `min_samples_leaf` parameter instead. This is what scikit-learn users do when pruning algorithms are unavailable. As discussed earlier, it is not easy to find a good value for this hyperparameter that performs similarly across a variety of datasets. In this case the result is very aggressive pruning except for datasets wine, iris and hepatitis. This is simply because the unpruned trees of these datasets were small to begin with, and this criterion aims for a low absolute number of nodes. The other configurations also score worse on these datasets because there is less potential for pruning if there are fewer nodes to start with. Table 6.1 gives an overview of the mean unpruned tree size of each dataset.

The previous heatmap only showed results obtained from the scikit-learn configu-

Dataset	Size
activity	580.94
diabetes	233.28
german_credit	329.32
heart-statlog	78.22
hepatitis	15.40
houses	499.02
ionosphere	37.06
iris	16.28
letter	3916.08
monks-problems-2	91.76
tic-tac-toe	143.10
vote	26.32
wdbc	34.78
wine	15.00

TABLE 6.1: Mean size of unpruned trees (i.e., configuration *none*) for each dataset.

rations. The next step is to compare these results with those of Weka. [Figure 6.2](#) visualizes the differences in number of nodes.

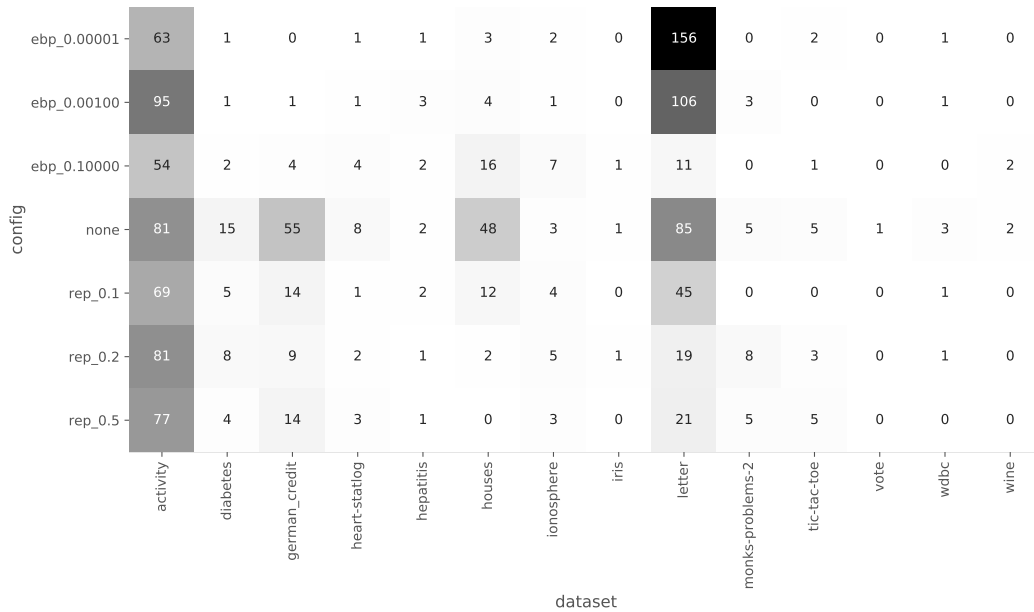


FIGURE 6.2: Absolute values of the differences between scikit-learn and Weka mean tree size per configuration and per dataset.

Most of the differences are minimal. The only configuration with some larger differences is the one where pruning is not applied. The previous heatmap showed

us that mean unpruned tree size can be a multiple of any mean pruned tree size. Consequently, the differences in size between unpruned trees are also larger in absolute numbers. Also note that this configuration does not involve any relevant custom code developed for this thesis. Additionally, it is proof that CART and C4.5 (and their respective implementations) do not produce the exact same trees, even when configured very similarly.

Column-wise, the letter, activity datasets — and to a lesser extent also houses and german_credit datasets — stand out. These are the datasets which on average generate the largest unpruned trees, as seen in [Table 6.1](#). Percentage-wise the effect pruning has on these four datasets is in line with the other datasets. As a result, their pruned trees are on average also large. This explains the larger absolute differences. When expressed relatively in percentages, the heatmap does not show any remarkable discrepancy between these four datasets and the rest.

In conclusion, the new pruning algorithms for scikit-learn decision trees behave as expected on the structure of the trees. The resulting sizes are grosso modo the same as when the experiments are performed in Weka. The next question is whether the performance also stays the same.

6.1.2 Accuracy

[Figure 6.3](#) shows a heatmap of the mean accuracy per scikit-learn configuration and per dataset. The most important observation here is that all configurations except for early stopping behave very similarly, despite large differences in the amount of pruning. Of course this is the purpose of pruning: reduce model complexity without affecting performance too much. This goal is clearly accomplished. At the same time, a strong contrast with the performance of the early stopping criterion is revealed. Across the board, it performs worse than all real pruning configurations. In the previous section we learned that early stopping acts on the tree structure the same way an aggressive pruning algorithm would. The main difference is that in general the pruning strategies are in a better position to decide whether a node is really useful. In the end, a tree of the same size but with a better node selection performs better than a tree that stopped growing prematurely because of some very generic criteria.

This heatmap also clearly shows that accuracy is strongly dependent on the dataset. Not all datasets can be modelled equally well with trees (or with any classifier for that matter). Differences in class balance also exist between datasets. Accuracy as a metric can be biased when the classes are strongly unbalanced. Alternative metrics such as the F1-score are better suited in such scenarios. A similar heatmap as shown here was generated using this F1-score, but it is not included because the results were almost identical.

A small note about the result of the *ebp_0.00001* configuration on the monks dataset: the low value for the confidence value hyperparameter makes the pruner

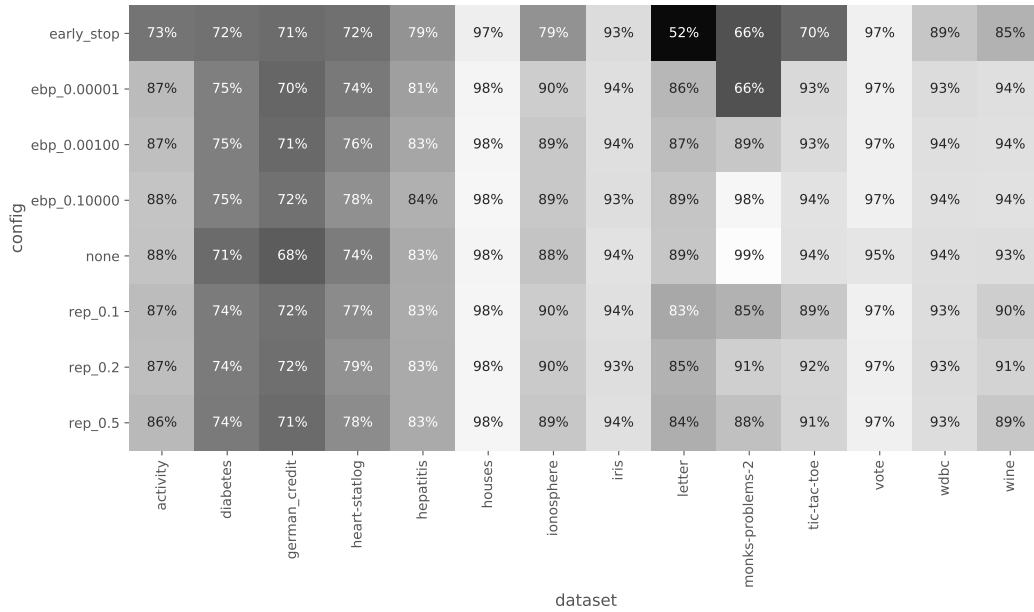


FIGURE 6.3: Mean accuracy per scikit-learn configuration and per dataset. Weka results are not included in this heatmap. The colour scale goes from 50% (black) to 100% (white).

very aggressive. In this case so much that the pruned tree consists of only one node that predicts the majority class all the time.

Figure 6.4 reveals that these scikit-learn accuracy scores are almost perfectly in line with the Weka scores. The only real concern is the eight percentage points difference in favour of scikit-learn on the monks dataset using Reduced Error Pruning with a 50% validation set size. Depending on which half of the data was withheld from training, the tree structure can vary dramatically. Consequently, the accuracy can vary a lot as well.

The same figure also reveals that the accuracy scores on the activity dataset are almost identical between scikit-learn and Weka configurations. Even the unpruned configurations perform almost identical. This contradicts earlier reports that noticed a significant difference in performance between the two tools. We cannot reproduce this issue.

6. RESULTS AND DISCUSSION

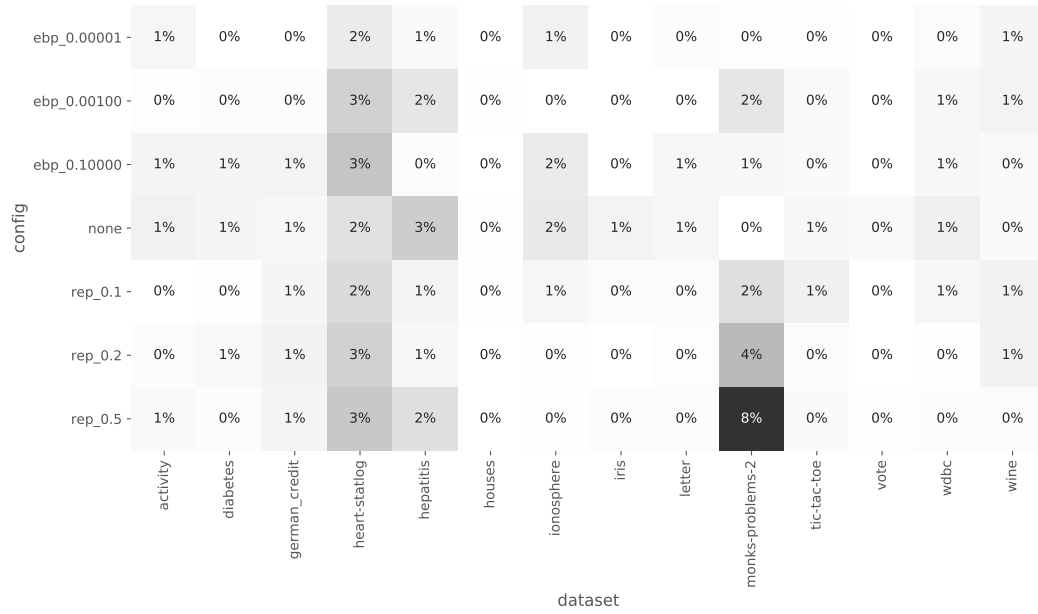


FIGURE 6.4: Absolute values of the differences between scikit-learn and Weka mean accuracy scores per configuration and per dataset. The colour scale goes from 0pp (white) to 10pp (black).

6.1.3 Fitting and training time

6.2 Categorical attributes

6.3 Binary versus non-binary trees

6.4 Conclusion

Chapter 7

Conclusion

Intro

7.1 Contributions

- overview of capabilities
- activity performance problem not reproduced
- pruning implementation (REP + EBP)
- CsvImporter: data pre-processor
- Categorical workaround performance test

7.2 Retrospective

7.3 Future work

- missing values fix inside DT
- cover more edge cases (multioutput, ...)

Appendices

Appendix A

The First Appendix

Appendices hold useful data which is not essential to understand the work done in the master's thesis. An example is a (program) source. An appendix can also have sections as well as figures and references.

Bibliography

- [BA97] Leonard A Breslow and David W Aha. Simplifying decision trees: A survey. *The Knowledge Engineering Review*, 12(1):1–40, 1997.
- [BFSO84] Leo Breiman, Jerome Friedman, Charles J Stone, and Richard A Olshen. *Classification and regression trees*. CRC press, 1984.
- [Bre01] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [CG16] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794. ACM, 2016.
- [Dem06] Janez Demsar. Statistical comparisons of classifiers over multiple data sets. 7:1–30, 01 2006.
- [EHWP16] F Eibe, MA Hall, IH Witten, and JC Pal. The weka workbench. *On-line appendix for "data mining: practical machine learning tools and techniques."*: Fourth Morgan Kaufmann, 2016.
- [EK01] Tapio Elomaa and Matti Kääriäinen. An analysis of reduced error pruning. *Journal of Artificial Intelligence Research*, 15:163–187, 2001.
- [Elo99] Tapio Elomaa. The biases of decision tree pruning strategies. In *International Symposium on Intelligent Data Analysis*, pages 63–74. Springer, 1999.
- [EMSK97] Floriana Esposito, Donato Malerba, Giovanni Semeraro, and J Kay. A comparative analysis of methods for pruning decision trees. *IEEE transactions on pattern analysis and machine intelligence*, 19(5):476–491, 1997.
- [ER99] Tapio Elomaa and Juho Rousu. General and efficient multisplitting of numerical attributes. *Machine learning*, 36(3):201–244, 1999.
- [Fog16] Agner Fog. The microarchitecture of intel, amd and via cpus. an optimization guide for assembly programmers and compiler makers. *Copenhagen University College of Engineering*, 2016.

- [GB18] Josh Gardner and Christopher Brooks. Evaluating predictive models of student success: Closing the methodological gap. *arXiv preprint arXiv:1801.08494*, 2018.
- [HMS66] Earl B Hunt, Janet Marin, and Philip J Stone. Experiments in induction. 1966.
- [HR76] Laurent Hyafil and Ronald L Rivest. Constructing optimal binary decision trees is np-complete. *Information processing letters*, 5(1):15–17, 1976.
- [KC01] Boonserm Kijssirikul and Kongsak Chongkasemwongse. Decision tree pruning using backpropagation neural networks. In *Neural Networks, 2001. Proceedings. IJCNN'01. International Joint Conference on*, volume 3, pages 1876–1880. IEEE, 2001.
- [KWM11] Jennifer R Kwapisz, Gary M Weiss, and Samuel A Moore. Activity recognition using cell phone accelerometers. *ACM SigKDD Explorations Newsletter*, 12(2):74–82, 2011.
- [KZP07] Sotiris B Kotsiantis, I Zaharakis, and P Pintelas. Supervised machine learning: A review of classification techniques. *Emerging artificial intelligence applications in computer engineering*, 160:3–24, 2007.
- [ld18] Scikit learn developers. Imputation of missing values. <http://scikit-learn.org/dev/modules/impute.html>, 2018. [Online; accessed 2018-05-22].
- [Man97] Yishay Mansour. Pessimistic decision tree pruning based on tree size. In *MACHINE LEARNING-INTERNATIONAL WORKSHOP THEN CONFERENCE-*, pages 195–201. Citeseer, 1997.
- [Min87] John Mingers. Rule induction with statistical data - a comparison with multiple regression. *Journal of the operational research Society*, 38(4):347–351, 1987.
- [Min89] John Mingers. An empirical comparison of pruning methods for decision tree induction. *Machine learning*, 4(2):227–243, 1989.
- [MRA⁺95] Manish Mehta, Jorma Rissanen, Rakesh Agrawal, et al. Mdl-based decision tree pruning. In *KDD*, volume 95, pages 216–221, 1995.
- [Mur98] Sreerama K Murthy. Automatic construction of decision trees from data: A multi-disciplinary survey. *Data mining and knowledge discovery*, 2(4):345–389, 1998.
- [NB86] T Niblett and I Bratko. Learning decision rules in noisy domains, 1986.
- [Oli06] Travis E Oliphant. *A guide to NumPy*, volume 1. Trelgol Publishing USA, 2006.

-
- [PVG⁺11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [QR89] J Ross Quinlan and Ronald L Rivest. Inferring decision trees using the minimum description length principle. *Information and computation*, 80(3):227–248, 1989.
- [Qui79] J Ross Quinlan. Discovering rules by induction from large collections of examples. *Expert systems in the micro electronics age*, 1979.
- [Qui83] J Ross Quinlan. Learning efficient classification procedures and their application to chess end games. In *Machine Learning, Volume I*, pages 463–482. Elsevier, 1983.
- [Qui86] J. Ross Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986.
- [Qui87] J. Ross Quinlan. Simplifying decision trees. *International journal of man-machine studies*, 27(3):221–234, 1987.
- [Qui93] J Ross Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.
- [Ris78] Jorma Rissanen. Modeling by shortest data description. *Automatica*, 14(5):465–471, 1978.
- [RM05] Lior Rokach and Oded Maimon. Top-down induction of decision trees classifiers-a survey. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 35(4):476–487, 2005.
- [Sha48] Claude E Shannon. A mathematical theory of communication (parts i and ii). *Bell System Tech. J.*, 27:379–423, 1948.
- [SSS⁺17] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.
- [VvRBT13] Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. Openml: Networked science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013.
- [ZXW⁺16] Matei Zaharia, Reynold S. Xin, Patrick Wendell, Tathagata Das, Michael Armbrust, Ankur Dave, Xiangrui Meng, Josh Rosen, Shrivaram Venkataraman, Michael J. Franklin, Ali Ghodsi, Joseph Gonzalez, Scott Shenker, and Ion Stoica. Apache spark: A unified engine for big data processing. *Commun. ACM*, 59(11):56–65, October 2016.

Master's thesis filing card

Student: Ir. Sven Van Hove

Title: Complete decision tree induction functionality in scikit-learn

Dutch title: Volledige beslissingsboom inductie functionaliteit in scikit-learn

UDC: 681.3*I20

Abstract:

500 word abstract

Thesis submitted for the degree of Master of Science in Artificial Intelligence, option Engineering and Computer Science

Thesis supervisors: Prof. dr. Jesse Davis

Prof. dr. ir. Hendrik Blockeel

Assessor: Dr. ir. Marc Claesen

Mentor: Elia Van Wolputte