

Applying the ID3 decision tree classifier to the mushrooms data set

Authors:

Yoei van Bruchem - s1013393

Mick van Hulst - s1013954

January 2018

Abstract

This project implemented the ID3 decision tree algorithm on a mushroom dataset. Using the generated tree, predictions were made regarding mushrooms being poisonous or edible. This algorithm was implemented by Yoei van Bruchem and Mick van Hulst. The algorithm was implemented in Python, using the course notes and description of Saed Sayad as a baseline. The tree is created by recursively using either a pre- or post-pruning approach (also called the training phase). After this phase, the test data set is classified using the generated tree. The algorithm classified the test data set perfectly using post-pruning.

1 Introduction

This report describes the algorithm that Yoei van Bruchem and Mick van Hulst have developed during their project for the data mining course at the Radboud University. The algorithm has been developed using the programming language Python.

The algorithm that was developed is called the ID3 decision tree algorithm. The algorithm is used for generating decision trees. In the case of this project, the focus lies on handling discrete variables only, meaning that the algorithm cannot handle continuous variables.

The code for this project can be found at the following link: https://github.com/mickvanhulst/dm_id3_dt_project. The README file in the GitHub repo explains the steps the user has to take to use this algorithm. This report will refer to the important parts in the code using the following notation: (file name, line numbers). Where an example could be: (main.py, 100-110). For this example the report would refer to the GitHub repository mentioned above, the file main.py and the line numbers 100 to 110.

2 Problem

The type of problems the ID3 decision tree algorithm tries to solve are supervised classification problems, meaning that if one has a set of training data for which the classes are known, one can train a decision tree with that data. Then when new data is gathered, classifications can be performed by using the previously generated decision tree.

3 Data set

During this project a mushroom data set has been used for the decision tree algorithm. The mushroom data set contains information of over 8,000 species of mushrooms. Each mushroom species is classified as either poisonous or edible and has various other characteristics that can be used in determining the edibility of the mushroom [2]. The algorithm will use the characteristics to classify records of mushrooms as either poisonous (p) or edible (e).

4 Related previous work

The related previous work that has been used can be split into two categories, where the first category of sources was used as background information to fully understand the theoretical/mathematical concepts behind the algorithm. This first category of previous work consists of the following sources:

- The lecture slides of the course were used as a baseline, meaning that it was used to find the corresponding formula's and the global process of the algorithm.
- The video's of Victor Lavrenko were used to learn more about the process of training and testing the decision tree [5].
- To implement post-pruning the web page of Saed Sayad was used. This created insight in the use of error estimation [4].

The second category consists of several previously developed projects that were used as example code. These projects were carefully analyzed and used to optimize the way of writing code for the algorithm. The projects are listed below:

- Jason Brownlee implemented a decision tree algorithm in Python using the Gini index. Although the Gini index wasn't implemented, the code gave insight in how to recursively call functions within Python [1].
- Gabriele Lazaro also implemented a decision tree in Python. In contrast to Jason Brownlee, she used the Information Gain. The code gave insight in the way of calculating the Information Gain [3].

5 Algorithm

The following subsections explain the ID3 decision tree algorithm and the process of splitting the data, training the decision tree and classifying the data set. These subsections also explain the mathematical concepts behind the algorithm and the variety of scenario's that may occur when running the algorithm.

5.1 User input

The first part of the algorithm consists of processing the inputs of the user. The algorithm requires the following inputs:

- The *data set* the user wants to use. The algorithm begins by splitting this data set into three groups. A training, test and validation set (ratio: 40/30/30). The training set is used to establish the tree, the validation set is used for post-pruning and the test set is used to test if the resulting model performs well or not (where 'well' is defined as the accuracy of classifying unknown records correctly or not).
- The *features and class label* of the data set. Using these values, the algorithm knows which features it can use to build the tree and which column (class label) it must use when comparing the predicted classes to the actual classes.

Using these inputs the algorithm proceeds to the training phase where it uses the features, data set and the class label to build the tree.

5.2 Training

During the training phase, the user can set several parameters which will determine how the decision tree will handle itself. The user can choose between post- and pre-pruning and can also decide if he/she wants to apply a maximum tree depth.

The algorithm starts by recursively calling the function that is responsible for generating the tree ('decision_tree.py', line 88-130). This recursive function uses the training data, the list of features and the class label as an input. Using this data, the algorithm determines the Information Gain for each feature at a particular split. The algorithm splits using the feature that has the maximum Information Gain.

$$GAIN_{split} = Entropy(p) - \left(\sum_{i=1}^k \frac{n_i}{n} Entropy(i) \right)$$

- $Entropy(p)$, the Entropy of the parent node.
- Parent node is split into k partitions.
- n_i is the number of records in partition i.

$$Entropy(t) = - \sum_c p(c|t) \log(p(c|t))$$

- $p(c|t)$ is the relative frequency of class c at node t .

After choosing which feature to split on, the same function is recursively called, using the remainder of the data and features yet to be split on as its inputs.

As mentioned above, the algorithm uses several inputs regarding pre-pruning, post-pruning and maximum tree depth. These inputs have an impact on how the algorithm handles itself during the training phase. The subsections below explain the impact of these inputs.

5.2.1 Pre-pruning

Pre-pruning is used to evaluate if a split is significant or not. This significance is determined by comparing the Information Gain of the previous split to the Information Gain of the to be made split. If the new Information Gain is lower or equal to the previous Information Gain, then the corresponding part of the tree gets pruned. In such a case, the algorithm will calculate the amount of classes at that particular split and return the class that occurs the most (i.e. the dominant class). This dominant class will be the leaf of the corresponding branch.

5.2.2 Max tree depth

The max tree depth is a predefined value which can be set by the user. This parameter decides how many layers of branches the tree can consist of. When the tree hits this corresponding tree depth, the tree will be pruned. In such a case the algorithm will calculate the amount of classes at that particular split and return the class that occurs the most (i.e. the dominant class). The most dominant class will be turned into a leaf.

5.2.3 Post-pruning

Post-pruning is a method that is used to avoid overfitting in previously build trees. Overfitting occurs when the decision tree algorithm fits the training data in a way that the tree would be inaccurate while predicting the class of the untrained data. Before post-pruning, the tree is first created completely, meaning that the tree only stops looking for newer features when the amount of classes left in the corresponding branch is equal to one [4].

Post-pruning consists of several steps:

1. Using the validation set, a validation tree is created. The validation set contains the original class. Using this, the algorithm can determine for each leaf if the instance of the validation set is classified correctly. This

step will result in a validated tree with a set of values consisting of the number of correctly and incorrectly classified instances.

2. Using the previously obtained information the algorithm applies the error estimation formula to estimate the error for each node ('post_pruning.py', line: 167-182). If the total error (calculated using the ratio of each of the leaves and it's corresponding error) under a node is greater than the error of that given node, then the algorithm makes the decision to cut that part of the tree off. The node will become a leaf with the most dominant class (of it's original sub leaves).

The error estimation formula consists of the following formula and corresponding components: $e = (f + \frac{z^2}{2N} + z\sqrt{\frac{f}{N} - \frac{f^2}{N} + \frac{z^2}{4N^2}})/(1 + \frac{z^2}{N})$

- f is the error on the training data.
- N is the number of instances covered by the leaf.
- z score from the normal distribution.

3. Now that the tree is pruned, the algorithm sanitizes the tree. Sanitation consists of removing all validation information and removes this with the corresponding classes. This step is necessary so that the remainder of the algorithm can recursively classify the test data.

5.3 Classification

Classification of the test data is also performed recursively. The corresponding part of the algorithm takes one row of the test data set and then recursively goes through all values of the corresponding features until it hits a leaf ('main.py', line: 132-159). If that's the case, then a class has been found.

After classifying all the test data, the predicted classes are compared to the actual classes, resulting in the accuracy.

6 Results

Using pre-pruning, the algorithm only needed one feature to classify around 98% of the test data correctly (see Figure 1). This feature is called 'odor'. Other than that, the post-pruning algorithm didn't prune anything, meaning that the algorithm didn't find a point in the tree which could be removed so that the error estimation would improve (see Figure 2). Post-pruning classified the test-set perfectly with an accuracy of 100%.

Figure 1: Result pre-pruning with 'odor' feature

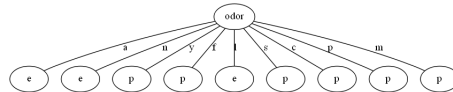
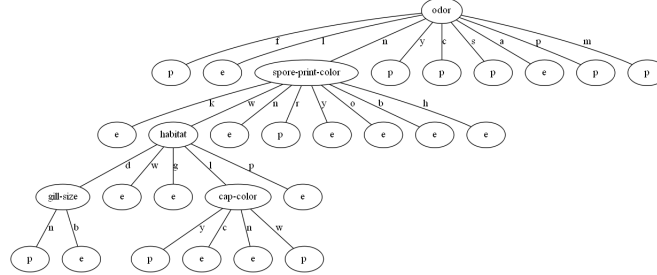
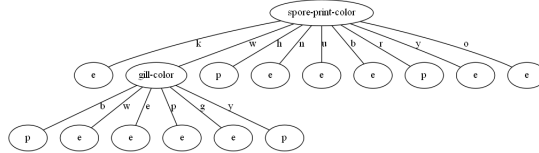


Figure 2: Result post-pruning with 'odor' feature



Another test was performed by removing the feature that seemed dominant (feature with name 'odor'). When using pre-pruning this resulted in a larger tree (see Figure 3), but in the decrease of the accuracy (around 94%). Using post-pruning it once again resulted in not removing any part of the tree (see Figure 4). Furthermore, the accuracy remained the same compared to the previous test (100%). To ensure that the post-pruning algorithm worked, the algorithm was tested with an extreme hardcoded error input. This resulted in an empty tree, proving that the algorithm does prune the tree when the error rate was significant.

Figure 3: Result pre-pruning without 'odor' feature



```

graph TD
    Root(spore-print-color) -- h --> Node1(gill-size)
    Root -- w --> Node2(gill-color)
    Root -- n --> Node3(gill-size)
    Root -- b --> Node4(gill-size)
    Root -- o --> Node5(population)
    Root -- t --> Node6(e)
    Root -- u --> Node7(e)
    Root -- y --> Node8(e)

    Node1 -- b --> Node1_1(p)
    Node1 -- n --> Node1_2(e)

    Node2 -- b --> Node2_1(p)
    Node2 -- w --> Node2_2(habitat)
    Node2 -- p --> Node2_3(e)
    Node2 -- g --> Node2_4(e)
    Node2 -- e --> Node2_5(e)
    Node2 -- y --> Node2_6(p)

    Node3 -- b --> Node3_1(e)
    Node3 -- n --> Node3_2(population)

    Node4 -- b --> Node4_1(e)
    Node4 -- n --> Node4_2(population)

    Node5 -- s --> Node5_1(p)
    Node5 -- v --> Node5_2(e)

    Node6 -- s --> Node6_1(p)
    Node6 -- v --> Node6_2(cap-color)
    Node6_2 -- e --> Node6_2_1(e)

    Node7 -- u --> Node7_1(cap-surface)
    Node7_1 -- y --> Node7_1_1(p)
    Node7_1 -- n --> Node7_1_2(p)
    Node7_1 -- e --> Node7_1_3(e)

    Node8 -- w --> Node8_1(cap-color)
    Node8_1 -- e --> Node8_1_1(e)

    Node3_2 -- s --> Node3_2_1(stalk-shape)
    Node3_2_1 -- e --> Node3_2_1_1(e)
    Node3_2_1 -- v --> Node3_2_1_2(p)

    Node4_2 -- s --> Node4_2_1(cap-surface)
    Node4_2_1 -- e --> Node4_2_1_1(e)
    Node4_2_1 -- v --> Node4_2_1_2(stalk-root)
    Node4_2_1_2 -- e --> Node4_2_1_2_1(p)
    Node4_2_1_2 -- b --> Node4_2_1_2_2(e)

    Node5_1 -- s --> Node5_1_1(stalk-root)
    Node5_1_1 -- e --> Node5_1_1_1(e)
    Node5_1_1 -- b --> Node5_1_1_2(p)

    Node6_2_1 -- s --> Node6_2_1_1(p)
    Node6_2_1 -- y --> Node6_2_1_2(p)
    Node6_2_1 -- f --> Node6_2_1_3(e)

    Node7_1_1 -- s --> Node7_1_1_1(p)
    Node7_1_1 -- y --> Node7_1_1_2(p)
    Node7_1_1 -- f --> Node7_1_1_3(e)

    Node7_1_2 -- s --> Node7_1_2_1(p)
    Node7_1_2 -- y --> Node7_1_2_2(p)
    Node7_1_2 -- f --> Node7_1_2_3(e)

    Node7_1_3 -- s --> Node7_1_3_1(p)
    Node7_1_3 -- y --> Node7_1_3_2(p)
    Node7_1_3 -- f --> Node7_1_3_3(e)

    Node7_1_4 -- s --> Node7_1_4_1(p)
    Node7_1_4 -- y --> Node7_1_4_2(p)
    Node7_1_4 -- f --> Node7_1_4_3(e)

    Node7_1_5 -- s --> Node7_1_5_1(p)
    Node7_1_5 -- y --> Node7_1_5_2(p)
    Node7_1_5 -- f --> Node7_1_5_3(e)

    Node7_1_6 -- s --> Node7_1_6_1(p)
    Node7_1_6 -- y --> Node7_1_6_2(p)
    Node7_1_6 -- f --> Node7_1_6_3(e)

    Node7_1_7 -- s --> Node7_1_7_1(p)
    Node7_1_7 -- y --> Node7_1_7_2(p)
    Node7_1_7 -- f --> Node7_1_7_3(e)

    Node7_1_8 -- s --> Node7_1_8_1(p)
    Node7_1_8 -- y --> Node7_1_8_2(p)
    Node7_1_8 -- f --> Node7_1_8_3(e)

    Node7_1_9 -- s --> Node7_1_9_1(p)
    Node7_1_9 -- y --> Node7_1_9_2(p)
    Node7_1_9 -- f --> Node7_1_9_3(e)

    Node7_1_10 -- s --> Node7_1_10_1(p)
    Node7_1_10 -- y --> Node7_1_10_2(p)
    Node7_1_10 -- f --> Node7_1_10_3(e)

    Node7_1_11 -- s --> Node7_1_11_1(p)
    Node7_1_11 -- y --> Node7_1_11_2(p)
    Node7_1_11 -- f --> Node7_1_11_3(e)

    Node7_1_12 -- s --> Node7_1_12_1(p)
    Node7_1_12 -- y --> Node7_1_12_2(p)
    Node7_1_12 -- f --> Node7_1_12_3(e)

    Node7_1_13 -- s --> Node7_1_13_1(p)
    Node7_1_13 -- y --> Node7_1_13_2(p)
    Node7_1_13 -- f --> Node7_1_13_3(e)

    Node7_1_14 -- s --> Node7_1_14_1(p)
    Node7_1_14 -- y --> Node7_1_14_2(p)
    Node7_1_14 -- f --> Node7_1_14_3(e)

    Node7_1_15 -- s --> Node7_1_15_1(p)
    Node7_1_15 -- y --> Node7_1_15_2(p)
    Node7_1_15 -- f --> Node7_1_15_3(e)

    Node7_1_16 -- s --> Node7_1_16_1(p)
    Node7_1_16 -- y --> Node7_1_16_2(p)
    Node7_1_16 -- f --> Node7_1_16_3(e)

    Node7_1_17 -- s --> Node7_1_17_1(p)
    Node7_1_17 -- y --> Node7_1_17_2(p)
    Node7_1_17 -- f --> Node7_1_17_3(e)

    Node7_1_18 -- s --> Node7_1_18_1(p)
    Node7_1_18 -- y --> Node7_1_18_2(p)
    Node7_1_18 -- f --> Node7_1_18_3(e)

    Node7_1_19 -- s --> Node7_1_19_1(p)
    Node7_1_19 -- y --> Node7_1_19_2(p)
    Node7_1_19 -- f --> Node7_1_19_3(e)

    Node7_1_20 -- s --> Node7_1_20_1(p)
    Node7_1_20 -- y --> Node7_1_20_2(p)
    Node7_1_20 -- f --> Node7_1_20_3(e)

    Node7_1_21 -- s --> Node7_1_21_1(p)
    Node7_1_21 -- y --> Node7_1_21_2(p)
    Node7_1_21 -- f --> Node7_1_21_3(e)

    Node7_1_22 -- s --> Node7_1_22_1(p)
    Node7_1_22 -- y --> Node7_1_22_2(p)
    Node7_1_22 -- f --> Node7_1_22_3(e)

    Node7_1_23 -- s --> Node7_1_23_1(p)
    Node7_1_23 -- y --> Node7_1_23_2(p)
    Node7_1_23 -- f --> Node7_1_23_3(e)

    Node7_1_24 -- s --> Node7_1_24_1(p)
    Node7_1_24 -- y --> Node7_1_24_2(p)
    Node7_1_24 -- f --> Node7_1_24_3(e)

    Node7_1_25 -- s --> Node7_1_25_1(p)
    Node7_1_25 -- y --> Node7_1_25_2(p)
    Node7_1_25 -- f --> Node7_1_25_3(e)

    Node7_1_26 -- s --> Node7_1_26_1(p)
    Node7_1_26 -- y --> Node7_1_26_2(p)
    Node7_1_26 -- f --> Node7_1_26_3(e)

    Node7_1_27 -- s --> Node7_1_27_1(p)
    Node7_1_27 -- y --> Node7_1_27_2(p)
    Node7_1_27 -- f --> Node7_1_27_3(e)

    Node7_1_28 -- s --> Node7_1_28_1(p)
    Node7_1_28 -- y --> Node7_1_28_2(p)
    Node7_1_28 -- f --> Node7_1_28_3(e)

    Node7_1_29 -- s --> Node7_1_29_1(p)
    Node7_1_29 -- y --> Node7_1_29_2(p)
    Node7_1_29 -- f --> Node7_1_29_3(e)

    Node7_1_30 -- s --> Node7_1_30_1(p)
    Node7_1_30 -- y --> Node7_1_30_2(p)
    Node7_1_30 -- f --> Node7_1_30_3(e)

    Node7_1_31 -- s --> Node7_1_31_1(p)
    Node7_1_31 -- y --> Node7_1_31_2(p)
    Node7_1_31 -- f --> Node7_1_31_3(e)

    Node7_1_32 -- s --> Node7_1_32_1(p)
    Node7_1_32 -- y --> Node7_1_32_2(p)
    Node7_1_32 -- f --> Node7_1_32_3(e)

    Node7_1_33 -- s --> Node7_1_33_1(p)
    Node7_1_33 -- y --> Node7_1_33_2(p)
    Node7_1_33 -- f --> Node7_1_33_3(e)

    Node7_1_34 -- s --> Node7_1_34_1(p)
    Node7_1_34 -- y --> Node7_1_34_2(p)
    Node7_1_34 -- f --> Node7_1_34_3(e)

    Node7_1_35 -- s --> Node7_1_35_1(p)
    Node7_1_35 -- y --> Node7_1_35_2(p)
    Node7_1_35 -- f --> Node7_1_35_3(e)

    Node7_1_36 -- s --> Node7_1_36_1(p)
    Node7_1_36 -- y --> Node7_1_36_2(p)
    Node7_1_36 -- f --> Node7_1_36_3(e)

    Node7_1_37 -- s --> Node7_1_37_1(p)
    Node7_1_37 -- y --> Node7_1_37_2(p)
    Node7_1_37 -- f --> Node7_1_37_3(e)

    Node7_1_38 -- s --> Node7_1_38_1(p)
    Node7_1_38 -- y --> Node7_1_38_2(p)
    Node7_1_38 -- f --> Node7_1_38_3(e)

    Node7_1_39 -- s --> Node7_1_39_1(p)
    Node7_1_39 -- y --> Node7_1_39_2(p)
    Node7_1_39 -- f --> Node7_1_39_3(e)

    Node7_1_40 -- s --> Node7_1_40_1(p)
    Node7_1_40 -- y --> Node7_1_40_2(p)
    Node7_1_40 -- f --> Node7_1_40_3(e)

    Node7_1_41 -- s --> Node7_1_41_1(p)
    Node7_1_41 -- y --> Node7_1_41_2(p)
    Node7_1_41 -- f --> Node7_1_4
```

- The feature 'odor' could be a dominant feature, meaning that this single feature significantly increases the accuracy of the algorithm.
- The data has a low amount of variety. This hypothesis was formed when the tree wasn't pruned at all when applying post-pruning. Normally, not removing any part of the tree would result in an overfitted tree. In this case the test set was classified perfectly, showing that the tree wasn't being overfit.

7 Analysis results

Overall, we're very satisfied with the result and feel like we've learned a lot throughout this process. Implementing an algorithm from scratch instead of using a package like Sci-Kit gave us a new understanding on the several concepts that we've used. It also resulted in us respecting the developers of such packages even more, as we realized that these packages make our lives a whole lot easier.

8 Improvements

After completing the development of the algorithm, the team of two students considered the following improvements:

- The algorithm could be further tested by using a different data set. As mentioned, it seemed that there is a low amount of variety in the data set. Testing it on different data sets would give room to test whether or not removing part of the tree when post-pruning has a positive impact on the overall accuracy.
- Implementation of features like boosting which could possibly increase the accuracy of the algorithm significantly.
- Extending this implementation to work with continuous variables. The data set that was used did not contain any continuous variables, so the feature wasn't implemented. For an algorithm to perform for any data set, this feature would have to be included.

References

- [1] Jason Brownlee. How to implement the decision tree algorithm from scratch in python, Aug 2017.
- [2] Holt. Mushroom analysis, Jun 2017.
- [3] Gabriele Lanaro. Implementing decision trees in python, Mar 2016.
- [4] Saed Sayad. Decision tree - classification.
- [5] victorlavrenko. Id3 algorithm: how it works, Jan 2014.