

Assignment 2: Decision Trees Algorithm

H. Blum, D. Cavezza, A. Paudice and M. Rohbeck
Machine Learning CO395
Imperial College London

November 8, 2015

1 Training

Decision trees are a conceptually simple, yet powerful model for solving classification problems. The ID3 algorithm is one of the simplest training algorithms for decision trees: it builds the tree by choosing at each step the attribute that maximizes the Information Gain, a measure that represents the discriminatory power of an attribute over the class values. The problem we are given consists in classifying facial expressions into 6 emotion classes via decision trees; each facial expression is represented as a feature vector of 45 Action Units (AU), where each feature is a binary value that indicates the presence or absence of each AU in the example. We decompose the problem in 6 binary decision problems, one corresponding to the recognition of one specific emotion. Therefore, we use ID3 to train 6 different trees, each specific to one emotion.

Our implementation of ID3 is provided in the files `train.m`, `DTTrain.m`, `choose_best_attr.m`, and `sample_entropy.m`, containing the functions with the same name.

The function `train` is the entry point for the training algorithm. It takes as inputs the training examples, the array of attribute ids (in our case, an array from 1 to 45), and an array of labels of the same length as the number of examples; such labels are numbers from 1 to 6, each one indicating one of the emotions in the original dataset. It returns the 6 desired trees.

Each tree is trained on a transformed version of the input data: at the i -th iteration of the for loop, the labels are made binary; the vector `binary_targets` contains a 1 in the positions corresponding to the i -th emotion, a 0 in all the other positions. The function `DTTrain` is called to train each tree on such transformed data.

```
1 function [ T ] = train( examples, attributes, labels )
2 %Train 6 binary trees on the training set (examples, labels) and return
3 %them into the tree array T.
4
5 [m, n] = size(examples);
6 binary_targets = zeros(m, 1);
7
8 for i = 1:6
9
10     %Switch to binary labels
11     index = find(labels == i);
12     binary_targets(index) = 1;
13
14     %Train a tree to learn an emotion
15     T(i) = DTTrain(examples, attributes, binary_targets);
16
17 %     DrawDecisionTree(T(i));
18
19     binary_targets = zeros(m, 1);
20
21 end
22
23 end
```

The function `DTTrain` is a recursive function implementing the ID3 algorithm. Every time a recursion occurs, a check is performed on the entropy of the current node: if it is not 0, the attribute that maximizes the Information Gain is chosen and associated to the current node, then the recursion occurs.

In the first call it is passed the whole training sample, the array 1:45 of attribute ids, and the binary labels.

```
1 function [ tree ] = DTTrain( examples, attributes, binary_targets )
2 %Train a decision tree on the dataset (examples, binary_target)
3 %according to the ID3 algorithm. The examples have the attributes listed in
```

```
4 %attributes
```

First, an empty tree is initialized.

```
6 %Initialize an empty tree
7 tree.op = [];
8
9 tree.kids = [];
10 tree.class = [];
```

Then a check on the input's entropy is performed. If the sample entropy in the labels is 0, the algorithm is in a pure node, that is all the current examples belong to the same class: in such case, the algorithm must stop, and the value returned by the current node is the class value held by the majority of the examples. Likewise, the algorithm must stop if there are no more attributes to decide about. In these cases, `tree.op` and `tree.kids` are not set, so as to mark the node as a leaf.

```
12 if (sample_entropy(binary_targets) == 0 || isempty(attributes))
13     % either pure targets or no more attributes to base decision on
14
15     tree.class = maj_value(binary_targets);
```

Else, the attribute that yields the highest Information Gain is chosen

```
17 else
18
19     tree.op = choose_best_attr(examples, attributes, binary_targets);
```

and the examples are split by the value they contain in the chosen attribute.

```
1     tree.kids = cell(1, 2);
2
3     % raise up the kids
4
5     for j = 1:2
6
7         %fprintf('#Examples: %d \n', length(examples));
8         child_index = find(examples(:, tree.op) == (j - 1));
9         child_examples = examples(child_index, :);
10        child_binary_targets = binary_targets(child_index);
```

At this point, the set of examples with a specific value of the chosen attribute may be empty. In that case, the corresponding kid node is made leaf and the returned class value is the value held by the majority of the nodes in the current node.

```
36         if (isempty(child_examples))
37             % we can't train this child, make a leaf with the majority
38             % value of all training data coming to the parent
39
40             tree.kids{j}.op = [];
41             tree.kids{j}.kids = [];
42             tree.kids{j}.class = maj_value(binary_targets);
```

Otherwise, the corresponding subtree is created. Note that in the recursive calls, the chosen attribute is deleted from the list of candidate attributes, since its value is fixed for the samples in each of the kids (0 for the left kid and 1 for the right).

```
44         else
45
46             % remove the used attribute from the list
47             index = find(attributes == tree.op);
48             attributes(index) = [];
49
50             % recursive training with this child
51             tree.kids{j} = DTTrain(child_examples, attributes, ...
52                                     child_binary_targets);
53
54         end
```

At line 19, the above functions calls `choose_best_attr`. This function computes the Information Gain of every available attribute and returns the id of the attribute that yields the highest Information Gain.

```

1 function [ best_attr ] = choose_best_attr( examples, attributes, binary_targets )
2 %chooseBestAtt computes and return in bestAtt the attribute in attributes
3 %which determine the maximum information gain for the set
4 %(examples, binary_targets). Return -1 if the set is pure because no split
5 %is required
6
7 %Compute the sample entropy
8 n_ex = length(binary_targets);           %Dataset size
9 E = sample_entropy(binary_targets); %Sample entropy
10
11 % initial values
12 best_gain = 0;
13 best_attr = attributes(1); % for now this is as good as anything
14
15 for i = 1:length(attributes)
16
17     % examples where attribute is 1
18     one_attr = binary_targets(examples(:, attributes(i)) == 1);
19     % examples where attribute is 0
20     zero_attr = binary_targets(examples(:, attributes(i)) == 0);
21
22     E_partition = length(one_attr)/n_ex * sample_entropy(one_attr) ...
23                 + length(zero_attr)/n_ex * sample_entropy(zero_attr);
24
25     gain = E - E_partition;
26
27     if (gain > best_gain)
28
29         best_gain = gain;
30         best_attr = attributes(i);
31
32     end
33
34 end
35
36 end

```

Finally, the files `sample_entropy.m` and `maj_value.m` simply contain the computation of the entropy and the most frequent label value in a sample respectively.

2 Decision making strategies

Since we have different trees, and each of these trees may commit mistakes in its classification, their outputs may be contradictory. Specifically, more than one tree may return 1 in response to the same example, or all trees may return 0. Therefore, we have devised different decision policies to resolve conflicts.

2.1 Random choice

One possible strategy is choosing randomly: if more than a tree returns 1, the label of the example is chosen randomly among the classes corresponding to such trees; if all the trees return 0, the label is picked randomly among all the classes. The implementation of this strategy is in the file `testTreesRandomChoice.m`.

```

1 function [ predictions ] = testTreesRandomChoice(T,x2)
2 % TestTreesRandomChoice performs a random selection among the emotions of
3 % all the trees that return a prediction of 1, or among all the emotions if
4 % all the trees return 0.
5
6 [m,n] = size(x2);
7 binary_predictions = zeros(m,6);
8 predictions = zeros(m,1);
9
10 for i = 1:m
11     for j = 1:6
12         binary_predictions(i,j) = predictionBinaryTree(T(j),x2(i,:));
13     end
14     % If all the predictions are 0, pick an emotion randomly
15     if(max(binary_predictions(i,:))==0)
16         predictions(i) = randi(6);
17     % Else pick an emotion randomly among the 1s
18     else

```

```

19         candidateEmotions = find(binary_predictions(i,:));
20         predictions(i) = candidateEmotions(randi(length(candidateEmotions)));
21     end
22
23 end
24
25
26 end

```

The function `predictionBinaryTree` called in line 12 is defined in the same file: it passes through the branches of the tree corresponding to the attribute values of the input example and returns the binary label corresponding to the reached leaf.

```

28 function [ class ] = predictionBinaryTree( tree, x )
29 %predictionBinaryTree walks along the tree to find the classification of
30 %the instance x
31
32 if (isempty(tree.op)) %Check if it is a leaf
33
34     class = tree.class;
35
36 else %Test the attribute
37
38     x_to_test = x(tree.op);
39
40     if (x_to_test == 0) %Follow the left branch
41
42         class = predictionBinaryTree(tree.kids{1}, x);
43
44     else %Follow the right branch
45
46         class = predictionBinaryTree(tree.kids{2}, x);
47
48     end
49
50 end
51
52 end

```

This strategy is one of the simplest to consider in resolving conflicts. It is chosen as our baseline strategy.

2.2 Score-based choice

2.3 Tree confidence strategy

3 Cross-validation

4 Performance

5 Questions

=====

6 Implementations

From the set of 6 Decision Trees, each tree will return a binary classification for a given data sample. Therefore, we need an Algorithm that decides on Basis of the output of the 6 trees, which class the given sample falls into. In fact, this algorithm does not even have to depend on the output of the classification trees. For example, one could think of an algorithm that simply outputs a randomly chosen classification and will have an accuracy of 17% for a balanced test set. However, we searched for algorithms with a better performance.

6.1 Stacking

As we already implemented the basis function to train a decision tree, we tried to approximate the decision function (the function which outputs a class for a given example) by another tree which is trained on basis of the output of all 6 binary classification trees for the training set. Therefore, given a sample, the decision algorithm works as follows:

1. give the sample to the 6 binary classification trees
2. give the output of all 6 trees to the decision function tree and output it's output

6.1.1 Statistics

The stacked tree method achieves a general accuracy of 72.4% on clean data and 58.8% in noisy data.

105	9	2	1	9	6
37	145	1	5	3	7
21	4	79	1	4	10
13	10	4	177	5	7
39	13	3	6	64	7
18	5	13	6	8	157

(a) confusion matrix for clean data

19	7	16	4	9	33
13	130	7	8	2	27
13	16	98	11	4	45
10	14	13	141	1	30
21	8	4	6	43	28
13	8	21	9	11	158

(c) confusion matrix for noisy data

Class	Precision	Recall	F_1
1	48.3%	79.6%	57.6%
2	78.5%	73.7%	75.5%
3	76.5%	64.9%	69.8%
4	90.4%	82.5%	86.0%
5	66.3%	47.6%	54.7%
6	85.8%	76.4%	79.7%

(b) performance for the different classes with clean data

Class	Precision	Recall	F_1
1	23.5%	21.6%	NaN%
2	70.2%	68.1%	68.8%
3	60.8%	51.1%	54.8%
4	78.1%	67.6%	72.0%
5	59.9%	38.1%	45.4%
6	49.1%	71.8%	57.7%

(d) performance for the different classes with noisy data

Figure 1: statistics of the stacked tree

6.2 Probabilistic Trees

In this implementation, we assign a score to each leaf of every tree between 0 and 1. The score is correlated to the number of training examples that can be classified by the given leaf with respect to the total number of training examples in every class.

$$\text{score} = \frac{\text{\# correctly classified training data by this leaf}}{\text{\# training data in this tree}}$$

Of all binary classification trees returning 1, the decision algorithm will pick the class corresponding to the tree with the highest decision score. If no tree returns 1, it will pick at random.

```

1 function [ predictions ] = decide_by_score(trees, testset)
2 % decide_by_score performs a classification based on the fraction of
3 % correctly classified training examples given by the tree leafs
4
5 [m,n] = size(testset);
6 predictions = zeros(m,1);
7
8 for i = 1:m
9     % test in all trees, find the one with the best score
10    best_score = 0;
11    predictions(i) = NaN;
12    for t = 1:6
13        [pred, score] = prediction_with_score(trees(t), testset(i,:));
14        if pred == 1
15            % the tree recognises this item as his class
16            if score > best_score
17                predictions(i) = t;
18                best_score = score;
19            end
20        end
21    end
22    % If all the predictions are 0, pick a random class
23    if(isnan(predictions(i)))
24        predictions(i) = randi(6);

```

```

25     end
26 end
27
28 end

```

6.2.1 Statistics

The described method achieves a general accuracy of 73.2% on clean data and 64.4% in noisy data.

88	15	7	4	11	7
15	145	5	7	14	12
7	5	83	2	5	17
4	8	6	188	5	5
15	19	2	10	73	13
3	9	14	11	4	166

(a) confusion matrix for clean data

21	11	20	8	19	9
16	132	13	15	6	5
14	13	108	22	11	19
8	13	10	161	10	7
15	10	8	7	59	11
13	9	14	10	10	164

(c) confusion matrix for noisy data

Class	Precision	Recall	F_1
1	66.3%	63.1%	63.4%
2	72.4%	75.1%	73.3%
3	67.1%	65.9%	66.2%
4	79.1%	84.0%	81.0%
5	60.5%	55.1%	57.5%
6	83.3%	81.1%	82.0%

(b) performance for the different classes with clean data

Class	Precision	Recall	F_1
1	25.5%	23.3%	NaN
2	71.4%	69.7%	69.9%
3	64.4%	56.7%	59.2%
4	71.6%	77.6%	74.0%
5	51.0%	52.9%	50.7%
6	76.2%	74.1%	74.6%

(d) performance for the different classes with noisy data

Figure 2: statistics on the probabilistic tree

7 Pruning

These Figures show the error rate on validation and training set as the size (i. e. the number of leaves) of the trees increase. For clean data the actual loss in accuracy is, as can be seen, negligible. With the optimal size, the error rate is 21% instead of 24% with no pruning.

For noisy data, this difference becomes much bigger. Pruning can increase performance in this case by approx. 10% from 62% to 68%.

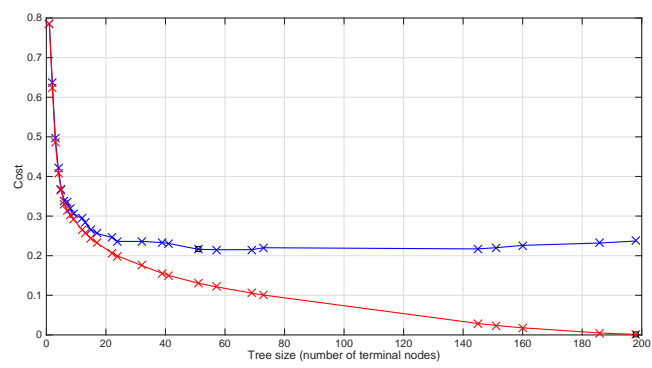
8 Performance with Noisy Data

The Performance with noisy data is much worse. Decision Trees can represent any kind of function, also the function generating noise. This is because while training the decision trees, we fit our tree to the training data, even if this data is noisy. Therefore noisy data generates noisy trees. For Class 1 (precision is less than half in comparison to clean data) we especially have a problem of overfitting which leads to much worse precision than with the clean data.

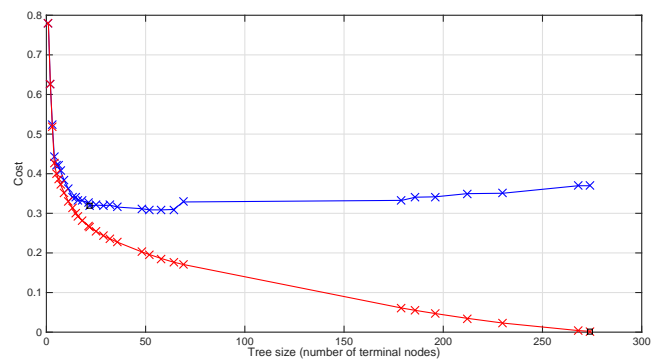
9 Choosing the Best Algorithm

We tested each of our algorithms using cross-validation. The Accuracy on the cross-validation test set is at this point a valid approximation of the accuracy on the unknown data set. On basis of these approximations, we choose the algorithm with the best accuracy. After choosing an algorithm, we train it with the whole data set we got. Unfortunately, we don't have a test set to estimate the performance of this trained algorithm on unseen data.

References



(a) Pruning Example with clean data



(b) Pruning Example with noisy data