# Machine Learning CS 529 Project 1 Report

## Esteban Guillen

## High-level Description of code

I implemented the ID3 algorithm in Python (Figure 1).  All code in contained in a file called ***id3.py***.  The code follows the algorithm ( https://en.wikipedia.org/wiki/ID3_algorithm ) closely and adds the ability to switch between entropy and misclassification error for the impurity measure.  Split stopping using the Chi-squared test was also added.
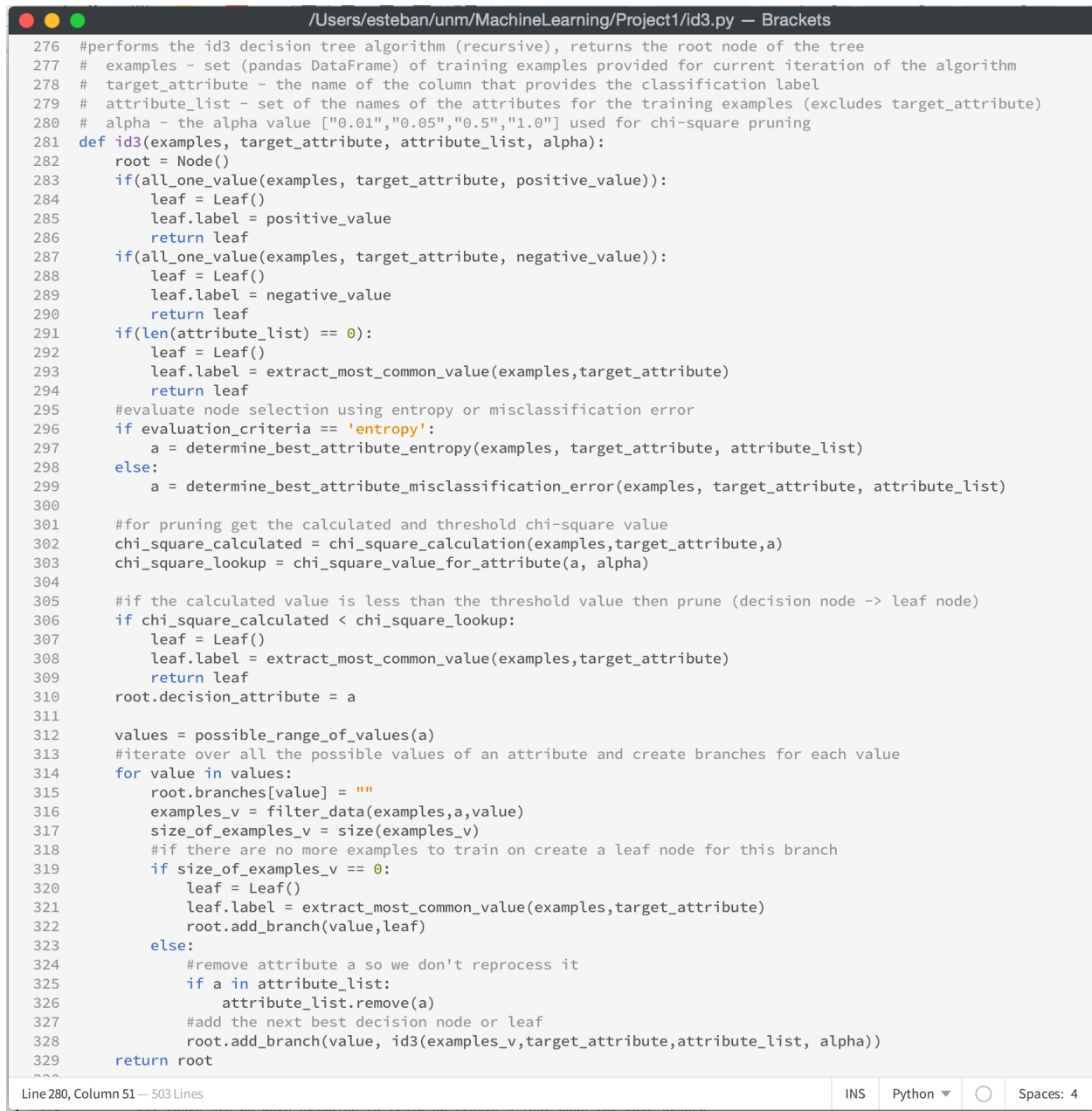
```python
276  #performs the id3 decision tree algorithm (recursive), returns the root node of the tree
277  #  examples - set (pandas DataFrame) of training examples provided for current iteration of the algorithm
278  #  target_attribute - the name of the column that provides the classification label
279  #  attribute_list - set of the names of the attributes for the training examples (excludes target_attribute)
280  #  alpha - the alpha value ["0.01","0.05","0.5","1.0"] used for chi-square pruning
281  def id3(examples, target_attribute, attribute_list, alpha):
282      root = Node()
283      if(all_one_value(examples, target_attribute, positive_value)):
284          leaf = Leaf()
285          leaf.label = positive_value
286          return leaf
287      if(all_one_value(examples, target_attribute, negative_value)):
288          leaf = Leaf()
289          leaf.label = negative_value
290          return leaf
291      if(len(attribute_list) == 0):
292          leaf = Leaf()
293          leaf.label = extract_most_common_value(examples,target_attribute)
294          return leaf
295      #evaluate node selection using entropy or misclassification error
296      if evaluation_criteria == 'entropy':
297          a = determine_best_attribute_entropy(examples, target_attribute, attribute_list)
298      else:
299          a = determine_best_attribute_misclassification_error(examples, target_attribute, attribute_list)
300
301      #for pruning get the calculated and threshold chi-square value
302      chi_square_calculated = chi_square_calculation(examples,target_attribute,a)
303      chi_square_lookup = chi_square_value_for_attribute(a, alpha)
304
305      #if the calculated value is less than the threshold value then prune (decision node -> leaf node)
306      if chi_square_calculated < chi_square_lookup:
307          leaf = Leaf()
308          leaf.label = extract_most_common_value(examples,target_attribute)
309          return leaf
310      root.decision_attribute = a
311
312      values = possible_range_of_values(a)
313      #iterate over all the possible values of an attribute and create branches for each value
314      for value in values:
315          root.branches[value] = ""
316          examples_v = filter_data(examples,a,value)
317          size_of_examples_v = size(examples_v)
318          #if there are no more examples to train on create a leaf node for this branch
319          if size_of_examples_v == 0:
320              leaf = Leaf()
321              leaf.label = extract_most_common_value(examples,target_attribute)
322              root.add_branch(value,leaf)
323          else:
324              #remove attribute a so we don't reprocess it
325              if a in attribute_list:
326                  attribute_list.remove(a)
327              #add the next best decision node or leaf
328              root.add_branch(value, id3(examples_v,target_attribute,attribute_list, alpha))
329      return root
```

Line 280, Column 51 — 503 Lines          INS     Python ▾     ◯     Spaces: 4

*Figure 1. ID3 Algorithm*

I used pandas DataFrame data structures for holding the training, testing and validation data (Figure 2 shows training and testing data being loaded). The DataFrame made it easy to index into the data and create subsets based on attribute values.

```python
 7    names_of_attributes = names=["label","cap-shape","cap-surface","cap-color","bruises","oder",
 8                                 "gill-attachment","gill-spacing","gill-size","gill-color","stalk-shape",
 9                                 "stalk-root","stalk-surface-above-ring","stalk-surface-below-ring",
10                                 "stalk-color-above-ring","stalk-color-below-ring","veil-type",
11                                 "veil-color","ring-number","ring-type","spore-print-color",
12                                 "population","habitat"]
13
14    #read training data into pandas DataFrame
15    training_data = pd.read_table("data/training.txt", sep=",", names=names_of_attributes)
16
17    #read testing data into pandas DataFrame
18    testing_data = pd.read_table("data/testing.txt", sep=",", names=names_of_attributes)
19
```

*Figure 2. Using pandas DataFrame*

I used two types of objects to represent my decision tree (Figure 3). **Node** which is a decision node that consists of a string representing the decision attribute and a dictionary to manage the branches and child nodes. **Leaf** which simply consists of a string that represents the classification label (e or p).

```python
21    #represents a decision node in the tree
22    class Node:
23        #identifies which attribute this node tests on
24        decision_attribute = ""
25
26        #initialize node
27        def __init__(self):
28            self.decision_attribute = ""
29            self.branches = {}
30
31        #adds a child Node or Leaf
32        def add_branch(self,attribute_value,node):
33            self.branches[attribute_value] = node
34
35        #returns all branches
36        def get_branches(self):
37            return self.branches
38
39    #represents a leaf node in the tree
40    class Leaf:
41        #label representing the classification value (e or p)
42        label = ""
```

*Figure 3. Objects used to represent the decision tree*

The calculations for entropy and misclassification error can be seen below (Figure 4). The calculations were very straight forward (entropy used the log function from the math library).

```python
178
179    #calculates the entropy on a subset (examples_v) of data
180    def calculate_entropy(examples_v,target_attribute):
181        entropy = 0.0
182        total_size = size(examples_v)
183        #if there is no data just return a zero value (does not contribute to calculation)
184        if total_size == 0:
185            return 0.0
186        positive_filtered_examples_v = filter_data(examples_v, target_attribute, positive_value)
187        negative_filtered_examples_v = filter_data(examples_v, target_attribute, negative_value)
188
189        num_positive = size(positive_filtered_examples_v)
190        calc_positive = 0.0
191        if num_positive != 0:
192            calc_positive = -(num_positive/total_size)*math.log((num_positive/total_size),2)
193
194
195        num_negative = size(negative_filtered_examples_v)
196        calc_negative = 0.0
197        if num_negative != 0:
198            calc_negative = - (num_negative/total_size)*math.log((num_negative/total_size),2)
199
200        #returns -p(+)log p(+) - p(-)log p(-)
201        return calc_positive + calc_negative
202
203    #calculates the misclassification error on a subset (examples_v) of data
204    def calculate_misclassification_error(examples_v, target_attribute):
205        error = 0.0
206        total_size = size(examples_v)
207        #if there is no data just return a zero value (does not contribute to calculation)
208        if total_size == 0:
209            return 0.0
210
211        positive_filtered_examples_v = filter_data(examples_v,target_attribute,positive_value)
212        negative_filtered_examples_v = filter_data(examples_v,target_attribute,negative_value)
213
214        num_positive = size(positive_filtered_examples_v)
215        p_positive = (num_positive/total_size)
216
217        num_negative = size(negative_filtered_examples_v)
218        p_negative = (num_negative/total_size)
219
220        error = 1.0 - max([p_positive,p_negative])
221        #returns 1 - max(probability of positive,probability of negative)
222        return error
```

Line 217, Column 54 — 503 Lines          INS    Python ▼  ○    Spaces: 4

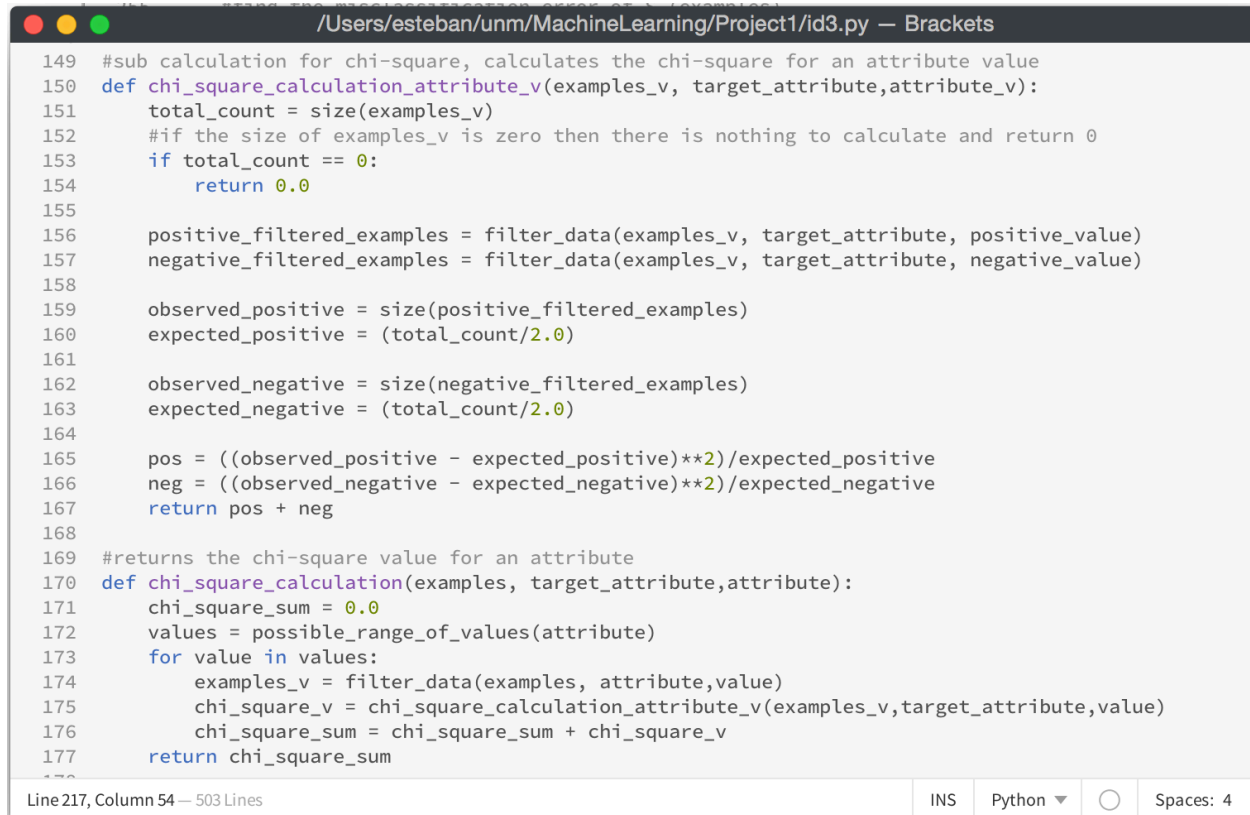*Figure 4. Entropy and misclassification calculations*

Gain was used to determine the "best attribute" for each recursive call of the ID3 algorithm. The implementations for entropy and misclassification can be seen below (Figure 5).

```python
224  #returns the attribute with the best informatin gain using entropy as the impurity measure
225  def determine_best_attribute_entropy(examples, target_attribute, attributes):
226      best_attribute = ""
227      max_information_gain = 0.0
228
229      #find the entropy of S (examples)
230      entropy_s = calculate_entropy(examples,target_attribute)
231      size_of_s = size(examples)
232
233
234      #loop through attributes and calculate the gain
235      for attribute in attributes:
236          values = possible_range_of_values(attribute)
237          information_gain = entropy_s
238          for value in values:
239              examples_v = filter_data(examples,attribute,value)
240              entropy_v = calculate_entropy(examples_v,target_attribute)
241              size_of_examples_v = size(examples_v)
242              weight = (size_of_examples_v/size_of_s)
243              information_gain = information_gain - weight*entropy_v
244          #if we have a new candidate for the best attribute store the gain and attribute name
245          if information_gain > max_information_gain:
246              max_information_gain = information_gain
247              best_attribute = attribute
248
249      return best_attribute
250
251  #returns the attribute with the best gain using misclassification error as the impurity measure
252  def determine_best_attribute_misclassification_error(examples, target_attribute, attributes):
253      best_attribute = ""
254      max_gain = 0.0
255      #find the misclassification error of S (examples)
256      error_s = calculate_misclassification_error(examples,target_attribute)
257      size_of_s = size(examples)
258
259      #loop through attributes and calculate the gain
260      for attribute in attributes:
261          values = possible_range_of_values(attribute)
262          gain = error_s
263          for value in values:
264              examples_v = filter_data(examples,attribute,value)
265              error_v = calculate_misclassification_error(examples_v,target_attribute)
266              size_of_examples_v = size(examples_v)
267              weight = (size_of_examples_v/size_of_s)
268              gain = gain - weight*error_v
269          #if we have a new candidate for the best attribute store the gain and attribute name
270          if gain > max_gain:
271              max_gain = gain
272              best_attribute = attribute
273
274      return best_attribute
```

*Figure 5. Gain calculations for entropy and misclassification*

The Chi-square test was used for split stopping and the code for calculating the Chi-square value is shown below (Figure 6). The calculation was broken up into 2 functions, the first calculates the Chi-square for a specific value of an attribute and the second sums all those calculations (over all possible values of an attribute).

```python
149   #sub calculation for chi-square, calculates the chi-square for an attribute value
150   def chi_square_calculation_attribute_v(examples_v, target_attribute,attribute_v):
151       total_count = size(examples_v)
152       #if the size of examples_v is zero then there is nothing to calculate and return 0
153       if total_count == 0:
154           return 0.0
155
156       positive_filtered_examples = filter_data(examples_v, target_attribute, positive_value)
157       negative_filtered_examples = filter_data(examples_v, target_attribute, negative_value)
158
159       observed_positive = size(positive_filtered_examples)
160       expected_positive = (total_count/2.0)
161
162       observed_negative = size(negative_filtered_examples)
163       expected_negative = (total_count/2.0)
164
165       pos = ((observed_positive - expected_positive)**2)/expected_positive
166       neg = ((observed_negative - expected_negative)**2)/expected_negative
167       return pos + neg
168
169   #returns the chi-square value for an attribute
170   def chi_square_calculation(examples, target_attribute,attribute):
171       chi_square_sum = 0.0
172       values = possible_range_of_values(attribute)
173       for value in values:
174           examples_v = filter_data(examples, attribute,value)
175           chi_square_v = chi_square_calculation_attribute_v(examples_v,target_attribute,value)
176           chi_square_sum = chi_square_sum + chi_square_v
177       return chi_square_sum
```

Line 217, Column 54 — 503 Lines        INS     Python ▾   ○    Spaces: 4

*Figure 6. Chi-square calculation code*

I implemented a number of helper functions to make the code more readable and maintainable. Check the full source code (***id3.py***) for more details.

**Accuracy Results**

All of my decision trees (entropy and misclassification error for 99, 95, 50, and 0 CL) produced 100% accuracy results (Figure 7). None of my trees got pruned using the Chi-squared test, all of my calculated values were less then the threshold values. The trees that were produced were small (4 levels for entropy and 5 levels for misclassification error, with most decision nodes having only one branch pointing to another decision node). The training data provided must have produced a near optimal tree, and no amount of pruning would improve performance (can't get better than 100%). Most of the training data was classified at the root node and 8 of the 9 branches pointed to Leaf nodes.

```
Estebans-MacBook-Pro:Project1 esteban$ python3 id3.py

Accuracy (using entropy and confidence level of 0):    100.0 %
   Number correctly classified:    2031
   Number incorrectly classified:  0
Accuracy (using entropy and confidence level of 99):   100.0 %
   Number correctly classified:    2031
   Number incorrectly classified:  0
Accuracy (using entropy and confidence level of 95):   100.0 %
   Number correctly classified:    2031
   Number incorrectly classified:  0
Accuracy (using entropy and confidence level of 50):    100.0 %
   Number correctly classified:    2031
   Number incorrectly classified:  0

Accuracy (using misclassification and confidence level of 0):    100.0 %
   Number correctly classified:    2031
   Number incorrectly classified:  0
Accuracy (using misclassification and confidence level of 99):   100.0 %
   Number correctly classified:    2031
   Number incorrectly classified:  0
Accuracy (using misclassification and confidence level of 95):   100.0 %
   Number correctly classified:    2031
   Number incorrectly classified:  0
Accuracy (using misclassification and confidence level of 50):    100.0 %
   Number correctly classified:    2031
   Number incorrectly classified:  0

best tree:  id3_entropy_confidence_level_0

classification predictions on validation data written to: validation-best-accuracy.txt

Estebans-MacBook-Pro:Project1 esteban$
```

*Figure 7. Classification results on the training.txt data*

I used the entropy with CL 0% decision tree to classify the validation data.  The 100% classification results can be seen below (Figure 8)
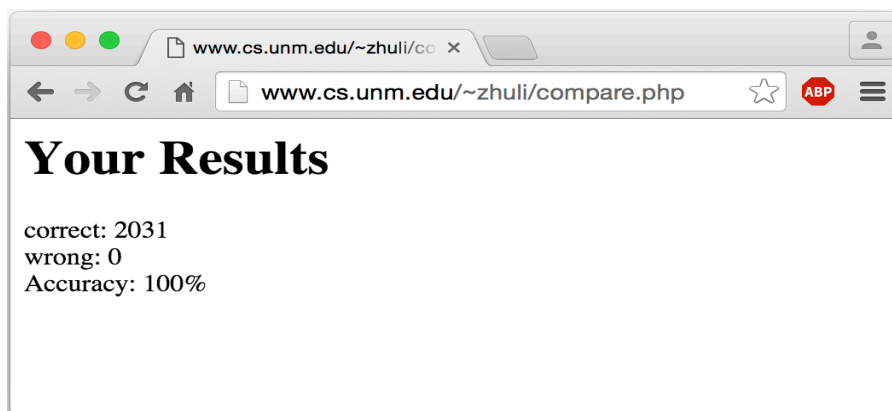


www.cs.unm.edu/~zhuli/compare.php

# Your Results

correct: 2031
wrong: 0
Accuracy: 100%

*Figure 8. Accuracy for the validation.txt data*

**Folklore Rules**

For proving or disproving the folklore rules I look at the data and/or my decision tree to come to a conclusion. I concluded that a folklore rule was true if the data or my decision tree backed-up the claim with a <u>high probability</u> (most of the data supported the claim). The included HTML export of my IPython Notebook shows how I queried the data.

<u>Rule 1</u>: Poisonous mushrooms are brightly colored:
**False**: I found there to be a fairly even split between edible and poisonous mushrooms for brightly colored (cap-color) mushrooms.

<u>Rule 2</u>: Poisonous mushrooms taste/smell bad:
**True**: The data provided strongly supports this (for smell bad). Almost all mushrooms with a bad odor were poisonous. My decision tree also supports this rule.

<u>Rule 3</u>: Poisonous mushrooms have a pointed or umbrella shaped cap:
**False**: The data provided strongly disproves this theory. Almost an even split between edible and poisonous mushroom having a pointed or umbrella shaped cap.

<u>Rule 4</u>: Edible mushroom have flat rounded shaped cap:
**False**: The data provided strongly disproves this theory. Almost an even split between edible and poisonous mushroom having a flat shaped cap.

<u>Rule 5</u>: Poisonous mushrooms have warts or scales on the cap:
**False**: The data provided strongly disproves this theory. Almost an even split between edible and poisonous mushroom having warts or scales on the cap.

<u>Rule 6</u>: Poisonous mushrooms have a bulbous cup or sac around the base:
**False**: The data provided strongly disproves this theory. Almost an even split between edible and poisonous mushroom having a bulbous cup or sac around the base.

<u>Rule 7</u>: Poisonous mushroom have a ring around the stem:
**False**: The data provided strongly disproves this theory. Almost an even split between edible and poisonous mushroom having a ring around the stem.

<u>Rule 8</u>: Poisonous mushrooms have gills that are thin and white:
**False**: The data provided strongly disproves this theory. Almost an even split between edible and poisonous mushroom having gills that are thin and white.

**Proposed Rules**

I could only find one rule not mentioned above to be true. The training data strongly (100%) supports "green spore print color mushroom are poisonous".
**Rule 9**: Poisonous mushrooms have green spore print color: