# ▾ Problem Statement

- A dataset collected in a cosmetics shop showing details of DEFINITION customers and whether or not they responded to a special offer to buy a new lip-stick is shown in table below.
- Use this dataset to build a decision tree, with Buys as the target variable, to help in buying lip-sticks in the future.
- Find the root node of decision tree.
- According to the decision tree you have made from previous training data set, what is the decision for the test data:
    - [Age < 21, Income = Low, Gender = Female, Marital Status = Married]?

| ID | Age | Income | Gender | Marital Status | Buys |
|----|------|--------|--------|----------------|------|
| 1  | <21  | High   | M      | Single         | N    |
| 2  | <21  | High   | M      | Married        | N    |
| 3  | 21-35| High   | M      | Single         | Y    |
| 4  | >35  | Medium | M      | Single         | Y    |
| 5  | >35  | Low    | F      | Single         | Y    |
| 6  | >35  | Low    | F      | Married        | N    |
| 7  | 21-35| Low    | F      | Married        | Y    |
| 8  | <21  | Medium | M      | Single         | N    |
| 9  | <21  | Low    | F      | Married        | Y    |
| 10 | >35  | Medium | F      | Single         | Y    |
| 11 | <21  | Medium | F      | Married        | Y    |
| 12 | 21-35| Medium | M      | Married        | Y    |
| 13 | 21-35| HIgh   | F      | Single         | Y    |
| 14 | >35  | Medium | M      | Married        | N    |

## Gini Impurity Formula

$$Gini = 1 - \sum_{i=1}^{C} (p_i)^2$$

## Information Gain Formula

- p = size_of_true_rows / (l size_of_true_rows + size_of_false_rows )
- current_uncertainty = gini impurity of input rows

- Compute the weighted average of the gini impurity of the true branch and the false branch, and subtract it from the current uncertainty

# information_gain = current_uncertainty - p * gini(true_rows) - (1 - p) * gini(false_rows)

## Decision Tree

- Full binary tree - each node contains either 2 or 0 child nodes
- Internal Nodes - Question
- Leaf Nodes - Final Prediction
- Edges - Answer to the Question

---

# CART Algorithm

---

## Algorithm: build_tree(rows)

- rows is the rows of the training_data passed to build_tree()
- find_best_split(rows) is used to determine the best question to ask
- Leaf is a class which returns final prediction
- partition(rows, question) is used to split the input rows based on the question asked
- Decision_Node is a class which contains the question, the true_branch and the false_branch

### Steps

1. information_gain, question = find_best_split(rows)
2. if information_gain := 0 then

    1. return Leaf(rows)
3. true_rows, false_rows = partition(rows, question)
4. true_branch = build_tree(true_rows)
5. false_branch = build_tree(false_rows)
6. return Decision_Node(question, true_branch, false_branch)

---

## Algorithm: find_best_split(rows)

- best_question is used to keep a track of the question with the highest information gain.
- best_gain is the information gain of the best_question.

- current_uncertainty is the gini impurity of the input rows.
- Question is a class which contains the input_feature and one of its unique values.
- partition is the function which splits the input rows into true_rows and false_rows based on the question asked.

## Steps

1. best_question = None

2. best_gain = 0

3. for each input_feature in input_features do

    1. unique_values = Get a set of unique values in the feature

    2. for each unique_value in unique_values

        3.2.1. question = Question(input_feature, unique_value)

        3.2.2. true_rows, false_rows := partition(rows, question)

        3.2.3. If question does not split the dataset then skip it

        3.2.4. information_gain := gain(true_rows, false_rows, current_uncertainty)

        3.2.5. if information_gain > best_gain then:

        3.2.5.1 best_gain := information_gain

        3.2.5.2 best_question := question

4. return best_question, best_gain

# ▾ Source Code

```
1 # Friend Functions
2
3 # Find the unique values for a column in a dataset.
4 def unique_vals(rows, col):
5     return set([row[col] for row in rows])
6
7
8 # Counts the number of each type of example in a dataset.
9 def class_counts(rows):
10     counts = {}  # a dictionary of label -> count.
11     for row in rows:
12         # in our dataset format, the label is always the last column
13         label = row[-1]
14         if label not in counts:
15             counts[label] = 0
```

```
16          counts[label] += 1
17      return counts
18
19 ##############################################################################
20
21 # A Question is used to partition a dataset.
22 # This class just records a 'column number' and a 'column value'.
23 # The 'match' method is used to compare the feature value
24 # in an example to the feature value stored in the question.
25 class Question:
26
27
28     def __init__(self, column, value):
29         self.column = column
30         self.value = value
31
32
33     # Compare the feature value in an example to the
34     # feature value in this question.
35     def match(self, example):
36         val = example[self.column]
37         return val == self.value
38
39
40     def __repr__(self):
41         # This is just a helper method to print
42         # the question in a readable format.
43         return "is {column} == {value}".format( column=header[self.column], value=
44
45 # A Leaf node classifies data.
46 # This holds a dictionary of class (e.g., "Apple") -> number of times
47 # it appears in the rows from the training data that reach this leaf.
48 class Leaf:
49
50     def __init__(self, rows):
51         self.predictions = class_counts(rows)
52
53 # A Decision Node asks a question.
54 # This holds a reference to the question, and to the two child nodes.
55 class Decision_Node:
56
57     def __init__(self,
58                  question,
59                  true_branch,
60                  false_branch):
61         self.question = question
62         self.true_branch = true_branch
63         self.false_branch = false_branch
64
65
66 ##############################################################################
67
```

```
 68  class Decision_Tree:
 69
 70
 71      # Partitions a dataset.
 72      # For each row in the dataset, check if it matches the question.
 73      # If so, add it to 'true rows', otherwise, add it to 'false rows'.
 74      def partition(self,rows, question):
 75          true_rows, false_rows = [], []
 76          for row in rows:
 77              if question.match(row):
 78                  true_rows.append(row)
 79              else:
 80                  false_rows.append(row)
 81          return true_rows, false_rows
 82
 83
 84      # Calculate the Gini Impurity for a list of rows.
 85      def gini(self,rows):
 86          counts = class_counts(rows)
 87          impurity = 1
 88          for lbl in counts:
 89              prob_of_lbl = counts[lbl] / float(len(rows))
 90              impurity -= prob_of_lbl**2
 91          return impurity
 92
 93
 94      # Information Gain. - The uncertainty of the starting node, minus the weighted
 95      def info_gain(self,left, right, current_uncertainty):
 96          p = float(len(left)) / (len(left) + len(right))
 97          return current_uncertainty - p * self.gini(left) - (1 - p) * self.gini(rig
 98
 99
100      # Find the best question to ask by iterating over every feature / value
101      # and calculating the information gain.
102      def find_best_split(self,rows):
103          best_gain = 0  # keep track of the best information gain
104          best_question = None  # keep train of the feature / value that produced it
105          current_uncertainty = self.gini(rows)
106          n_features = len(rows[0]) - 1  # number of columns
107
108          # For each feature
109          for col in range(n_features):
110
111              # unique values in the column
112              values = set([row[col] for row in rows])
113
114              # for each value
115              for val in values:
116
117                  # Ask Question
118                  question = Question(col, val)
119
```

```
119
120                    # Split the dataset
121                    true_rows, false_rows = self.partition(rows, question)
122
123                    # Skip this split if it doesn't divide the dataset.
124                    if len(true_rows) == 0 or len(false_rows) == 0:
125                        continue
126
127                    # Calculate the information gain from this split
128                    gain = self.info_gain(true_rows, false_rows, current_uncertainty)
129
130                    # Record the best gain and the best question
131                    if gain >= best_gain:
132                        best_gain, best_question = gain, question
133
134        # Return Question with the highest information gain
135        return best_gain, best_question
136
137
138    # Builds the tree.
139    # Rules of recursion:
140    #    1) Assume that it works.
141    #    2) Start by checking for the base case (no further information gain).
142    #    3) Prepare for giant stack traces.
143    def build_tree(self,rows):
144
145        # Try partitioing the dataset on each of the unique attribute,
146        # calculate the information gain,
147        # and return the question that produces the highest gain.
148        gain, question = self.find_best_split(rows)
149
150        # Base case: no further info gain
151        # Since we can ask no further questions,
152        # we'll return a leaf.
153        if gain == 0:
154            return Leaf(rows)
155
156        # If we reach here, we have found a useful feature / value
157        # to partition on.
158        true_rows, false_rows = self.partition(rows, question)
159
160        # Recursively build the true branch.
161        true_branch = self.build_tree(true_rows)
162
163        # Recursively build the false branch.
164        false_branch = self.build_tree(false_rows)
165
166        # Return a Question node.
167        # This records the best feature / value to ask at this point,
168        # as well as the branches to follow
169        # dependingo on the answer.
170        return Decision_Node(question, true_branch, false_branch)
```

```
171
172
173     def print_tree(self,node, spacing=""):
174
175         # Base case: we've reached a leaf
176         if isinstance(node, Leaf):
177             print (spacing + "Predict", node.predictions)
178             return
179
180         # Print the question at this node
181         print (spacing + str(node.question))
182
183         # Call this function recursively on the true branch
184         print (spacing + '--> True:')
185         self.print_tree(node.true_branch, spacing + "  ")
186
187         # Call this function recursively on the false branch
188         print (spacing + '--> False:')
189         self.print_tree(node.false_branch, spacing + "  ")
190
191
192     # Rules of recursion:
193     #   1) Assume that it works.
194     #   2) Start by checking for the base case (no further information gain).
195     #   3) Prepare for giant stack traces.
196     def classify(self,row, node):
197
198         # Base case: we've reached a leaf
199         if isinstance(node, Leaf):
200             return node.predictions
201
202         # Decide whether to follow the true-branch or the false-branch.
203         # Compare the feature / value stored in the node,
204         # to the example we're considering.
205         if node.question.match(row):
206             return self.classify(row, node.true_branch)
207         else:
208             return self.classify(row, node.false_branch)
209
210
211     def print_leaf(self,counts):
212         total = sum(counts.values()) * 1.0
213         probs = {}
214         for lbl in counts.keys():
215             probs[lbl] = str(int(counts[lbl] / total * 100)) + "%"
216         return probs
217
218     def set_root(self,node):
219       self.root = node
220
221     def get_root(self):
222       return self.root
```

```
222        return self.root
223
224  ############################################################################
225  class PCAG:
226
227
228    def __init__(self,training_data,testing_data,header):
229      training_data = self.data_validation(training_data,header)
230      DTC = self.training(training_data)
231      self.prediction(DTC, testing_data)
232
233
234    def data_validation(self,training_data,header):
235
236      print('\nData Validation\n===============')
237
238      print('# Raw Train Data = {length}'.format(length=len(training_data)))
239      training_data_cleaned = []
240      for record in training_data:
241        age, income, gender, marital_status, buys = record
242        if age not in ['<21','21-35','>35']:
243          pass
244        if income not in ['Low','Medium','High']:
245          pass
246        if gender not in ['M','F']:
247          pass
248        if marital_status not in ['Single','Maried']:
249          pass
250        if buys not in ['N','Y']:
251          pass
252        training_data_cleaned.append(
253            [age, income, gender, marital_status, buys]
254        )
255      print('# Cleaned Train Data = {length}'.format(length=len(training_data_cleane
256
257      return training_data_cleaned
258
259
260    def training(self,training_data):
261      print('\nTraining\n===============')
262      DTC = Decision_Tree()
263      decision_tree = DTC.build_tree(training_data)
264      DTC.set_root(decision_tree)
265      print('Final Tree')
266      DTC.print_tree(DTC.get_root())
267      return DTC
268
269
270    def prediction(self,DTC,testing_data):
271      print('\nPrediction\n===============')
272      print('Prediction')
273      for row in testing_data:
```

```
274        print("Record: {record}".format(record=row))
275        print("Predicted: {predicted}".format(predicted=DTC.print_leaf(DTC.classify(
276      pass
277
278 ############################################################################
279
280 training_data = [
281    ['<21', 'High', 'M', 'Single', 'N'],
282    ['<21', 'High', 'M', 'Married', 'N'],
283    ['21-35', 'High', 'M', 'Single', 'Y'],
284    ['>35', 'Medium', 'M', 'Single', 'Y'],
285    ['>35', 'Low', 'F', 'Single', 'Y'],
286    ['>35', 'Low', 'F', 'Married', 'N'],
287    ['21-35', 'Low', 'F', 'Married', 'Y'],
288    ['<21', 'Medium', 'M', 'Single', 'N'],
289    ['<21', 'Low', 'F', 'Married', 'Y'],
290    ['>35', 'Medium', 'F', 'Single', 'Y'],
291    ['<21', 'Medium', 'F', 'Married', 'Y'],
292    ['21-35', 'Medium', 'M', 'Married', 'Y'],
293    ['21-35', 'High', 'F', 'Single', 'Y'],
294    ['>35', 'Medium', 'M', 'Married', 'N']
295 ]
296 testing_data = [
297    ['<21', 'Low', 'F', 'Married']
298 ]
299 header = ["Age", "Income","Gender","Marital Status","Buys"]
300 pcag = PCAG(training_data,testing_data,header)
```

```
    Data Validation
    ================
    # Raw Train Data = 14
    # Cleaned Train Data = 14

    Training
    ================
    Final Tree
    is Age == 21-35
    --> True:
      Predict {'Y': 4}
    --> False:
      is Gender == F
      --> True:
        is Marital Status == Married
        --> True:
          is Age == >35
          --> True:
            Predict {'N': 1}
          --> False:
            Predict {'Y': 2}
        --> False:
          Predict {'Y': 2}
      --> False:
        is Age == >35
```

```
        --> True:
          is Marital Status == Married
          --> True:
            Predict {'N': 1}
          --> False:
            Predict {'Y': 1}
        --> False:
          Predict {'N': 3}


    Prediction
    ================
    Prediction
    Record: ['<21', 'Low', 'F', 'Married']
    Predicted: {'Y': '100%'}
```