

Example: Building a decision tree using Hunt's Algorithm

Impurity measure: Gini index

Splitting approach: Binary split

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Normal attributes class attribute

$$Gini(t) = 1 - \sum_j [p(j|t)]^2$$

Step 1 Determining the best split for the root node

① Refund = {Yes, No}

Refund

Yes No

N1 N2

Cheat → Yes = 0

Cheat → Yes = 3

Cheat → No = 3

Cheat → No = 4

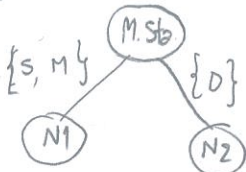
$$GINI(N1) = 1 - \left[\left(\frac{0}{3} \right)^2 + \left(\frac{3}{3} \right)^2 \right] = 0$$

$$GINI(N2) = 1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right] = 0.49$$

$$GINI_{split} = \left(\frac{3}{10} \right)(0) + \left(\frac{7}{10} \right)(0.49) = 0.343$$

② Marital Status = {Single, Married, Divorced}

2.1



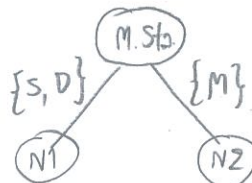
Cheat → Yes = 2 Cheat → Yes = 1
Cheat → No = 6 Cheat → No = 1

$$GINI(N1) = 1 - \left[\left(\frac{2}{8} \right)^2 + \left(\frac{6}{8} \right)^2 \right] = 0.375$$

$$GINI(N2) = 1 - \left[\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right] = 0.5$$

$$GINI_{split} = \left(\frac{8}{10} \right)(0.375) + \left(\frac{2}{10} \right)(0.5) = 0.4$$

2.2



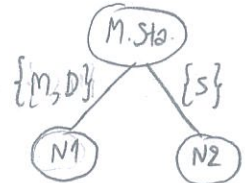
Cheat → Yes = 3 Cheat → Yes = 0
Cheat → No = 3 Cheat → No = 4

$$GINI(N1) = 1 - \left[\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right] = 0.5$$

$$GINI(N2) = 1 - \left[\left(\frac{0}{4} \right)^2 + \left(\frac{4}{4} \right)^2 \right] = 0$$

$$GINI_{split} = \left(\frac{6}{10} \right)(0.5) + \left(\frac{4}{10} \right)(0) = 0.3$$

2.3



Cheat → Yes = 1 Cheat → Yes = 2
Cheat → No = 5 Cheat → No = 2

$$GINI(N1) = 1 - \left[\left(\frac{1}{6} \right)^2 + \left(\frac{5}{6} \right)^2 \right] = 0.28$$

$$GINI(N2) = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

$$GINI_{split} = \left(\frac{6}{10} \right)(0.28) + \left(\frac{4}{10} \right)(0.5) = 0.368$$

Cheat \rightarrow N N N Y Y Y N N N
 (3) Taxable Income = {60K, 70K, 75K, 85K, 90K, 95K, 100K, 120K, 125K, 220K} (2)
 In this case, we have 10 distinct values, thus we have to use these values as a splitting point.

(3.1)

Cheat \rightarrow Yes = 0 Cheat \rightarrow Yes = 3
 Cheat \rightarrow No = 0 Cheat \rightarrow No = 7
 GINI(N1) = $1 - \left[\left(\frac{0}{0} \right)^2 + \left(\frac{0}{0} \right)^2 \right] = 1$
 GINI(N2) = $1 - \left[\left(\frac{3}{10} \right)^2 + \left(\frac{7}{10} \right)^2 \right] = 0.42$
 GINI_{split} = $\left(\frac{0}{10} \right)(1) + \left(\frac{10}{10} \right)(0.42) = 0.42$

(3.2)

Cheat \rightarrow Yes = 0 Cheat \rightarrow Yes = 3
 Cheat \rightarrow No = 1 Cheat \rightarrow No = 6
 GINI(N1) = $1 - \left[\left(\frac{0}{1} \right)^2 + \left(\frac{1}{1} \right)^2 \right] = 0$
 GINI(N2) = $1 - \left[\left(\frac{3}{9} \right)^2 + \left(\frac{6}{9} \right)^2 \right] = 0.44$
 GINI_{split} = $\left(\frac{1}{10} \right)(0) + \left(\frac{9}{10} \right)(0.44) = 0.4$

(3.3)

Cheat \rightarrow Yes = 0 Cheat \rightarrow Yes = 3
 Cheat \rightarrow No = 2 Cheat \rightarrow No = 5
 GINI(N1) = $1 - \left[\left(\frac{0}{2} \right)^2 + \left(\frac{2}{2} \right)^2 \right] = 0$
 GINI(N2) = $1 - \left[\left(\frac{3}{8} \right)^2 + \left(\frac{5}{8} \right)^2 \right] = 0.469$
 GINI_{split} = $\left(\frac{2}{10} \right)(0) + \left(\frac{8}{10} \right)(0.469) = 0.375$

(3.4)

Cheat \rightarrow Yes = 0 Cheat \rightarrow Yes = 3
 Cheat \rightarrow No = 3 Cheat \rightarrow No = 4
 GINI(N1) = $1 - \left[\left(\frac{0}{3} \right)^2 + \left(\frac{3}{3} \right)^2 \right] = 0$
 GINI(N2) = $1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right] = 0.49$
 GINI_{split} = $\left(\frac{3}{10} \right)(0) + \left(\frac{7}{10} \right)(0.49) = 0.343$

(3.5)

Cheat \rightarrow Yes = 1 Cheat \rightarrow Yes = 2
 Cheat \rightarrow No = 3 Cheat \rightarrow No = 4
 GINI(N1) = $1 - \left[\left(\frac{1}{4} \right)^2 + \left(\frac{3}{4} \right)^2 \right] = 0.375$
 GINI(N2) = $1 - \left[\left(\frac{2}{6} \right)^2 + \left(\frac{4}{6} \right)^2 \right] = 0.44$
 GINI_{split} = $\left(\frac{4}{10} \right)(0.375) + \left(\frac{6}{10} \right)(0.44) = 0.414$

(3.6)

Cheat \rightarrow Yes = 2 Cheat \rightarrow Yes = 1
 Cheat \rightarrow No = 3 Cheat \rightarrow No = 4
 GINI(N1) = $1 - \left[\left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right] = 0.48$
 GINI(N2) = $1 - \left[\left(\frac{1}{5} \right)^2 + \left(\frac{4}{5} \right)^2 \right] = 0.32$
 GINI_{split} = $\left(\frac{5}{10} \right)(0.48) + \left(\frac{5}{10} \right)(0.32) = 0.4$

(3.7)

Cheat \rightarrow Yes = 3 Cheat \rightarrow Yes = 0
 Cheat \rightarrow No = 3 Cheat \rightarrow No = 4
 GINI(N1) = $1 - \left[\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right] = 0.5$
 GINI(N2) = $1 - \left[\left(\frac{0}{4} \right)^2 + \left(\frac{4}{4} \right)^2 \right] = 0$
 GINI_{split} = $\left(\frac{6}{10} \right)(0.5) + \left(\frac{4}{10} \right)(0) = 0.3$

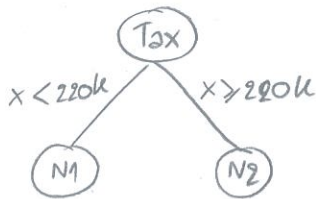
(3.8)

Cheat \rightarrow Yes = 3 Cheat \rightarrow Yes = 0
 Cheat \rightarrow No = 4 Cheat \rightarrow No = 3
 GINI(N1) = $1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right] = 0.49$
 GINI(N2) = $1 - \left[\left(\frac{0}{3} \right)^2 + \left(\frac{3}{3} \right)^2 \right] = 0$
 GINI_{split} = $\left(\frac{7}{10} \right)(0.49) + \left(\frac{3}{10} \right)(0) = 0.343$

(3.9)

Cheat \rightarrow Yes = 3 Cheat \rightarrow Yes = 0
 Cheat \rightarrow No = 5 Cheat \rightarrow No = 2
 GINI(N1) = $1 - \left[\left(\frac{3}{8} \right)^2 + \left(\frac{5}{8} \right)^2 \right] = 0.469$
 GINI(N2) = $1 - \left[\left(\frac{0}{2} \right)^2 + \left(\frac{2}{2} \right)^2 \right] = 0$
 GINI_{split} = $\left(\frac{8}{10} \right)(0.469) + \left(\frac{2}{10} \right)(0) = 0.375$

3.10



Cheat → Yes = 3 Cheat → Yes = 0
 Cheat → No = 6 Cheat → No = 1

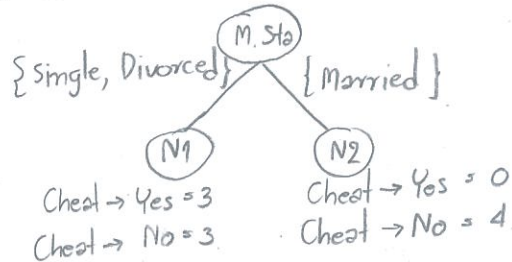
$$GINI(N1) = 1 - \left[\left(\frac{3}{9} \right)^2 + \left(\frac{6}{9} \right)^2 \right] = 0.44$$

$$GINI(N2) = 1 - \left[\left(\frac{0}{1} \right)^2 + \left(\frac{1}{1} \right)^2 \right] = 0$$

$$GINI_{split} = \left(\frac{9}{10} \right) (0.44) + \left(\frac{1}{10} \right) (0) = 0.4$$

Minimum { 0.343, 0.4, 0.3, 0.368, 0.42, 0.4, 0.375, 0.3 }
 { 0.343, 0.414, 0.4, 0.3, 0.343, 0.375, 0.4 } = 0.3

Since the minimum $GINI_{split}$ is obtained from test condition 2.2 and 3.7, we can use either test condition 2.2 or 3.7 as a test condition for the root node. In this case, I choose 2.2.



Step 2: Check whether we should stop splitting the records at node N1.

Two main stopping criteria:

- 1) Stop expanding a node when all the records belong to the same class
- 2) Stop expanding a node when all the records have identical attribute values.

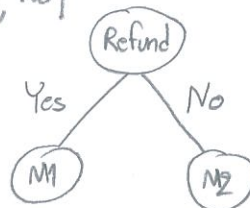
∴ Since all the records at node N1 do not satisfy the two main stopping criteria, then we need to split these records further.

Step 3: Check whether we should stop splitting the records at node N2.

∴ Since all the records at node N2 satisfy the first main stopping criterion, then we stop splitting and declare this node as a leaf node of value No.

Step 4: Determining the best split for node N1. (Now, we have only 6 records left at node N1)

① Refund = {Yes, No}



Cheat → Yes = 0 Cheat → Yes = 3
 Cheat → No = 2 Cheat → No = 1

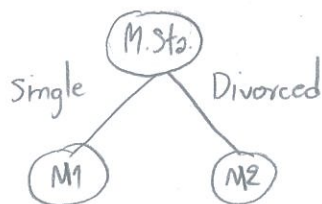
$$GINI(N1) = 1 - \left[\left(\frac{0}{2} \right)^2 + \left(\frac{2}{2} \right)^2 \right] = 0$$

$$GINI(N2) = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 0.375$$

$$GINI_{split} = \left(\frac{2}{6} \right) (0) + \left(\frac{4}{6} \right) (0.375) = 0.25$$

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

② Marital Status = { Single, Divorced }



Cheat → Yes = 2 Cheat → Yes = 1

Cheat → No = 2 Cheat → No = 1

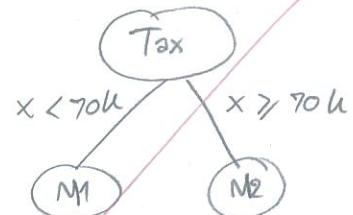
$$GINI(M1) = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

$$GINI(M2) = 1 - \left[\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right] = 0.5$$

$$GINI_{split} = \left(\frac{4}{6} \right) (0.5) + \left(\frac{2}{6} \right) (0.5) = 0.5$$

③ Taxable Income = { 70k, 85k, 90k, 95k, 125k, 220k } ④

3.1



Cheat → Yes = 0

Cheat → Yes = 3

Cheat → No = 0

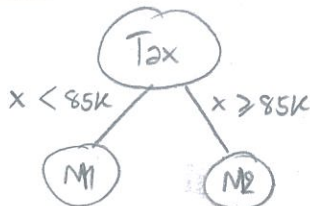
Cheat → No = 3

$$GINI(M1) = 1 - \left[\left(\frac{0}{0} \right)^2 + \left(\frac{0}{0} \right)^2 \right] \text{ (undefined) } \approx 1$$

$$GINI(M2) = 1 - \left[\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right] = 0.5$$

$$GINI_{split} = \left(\frac{0}{6} \right) (1) + \left(\frac{6}{6} \right) (0.5) = 0.5$$

3.2



Cheat → Yes = 0 Cheat → Yes = 3

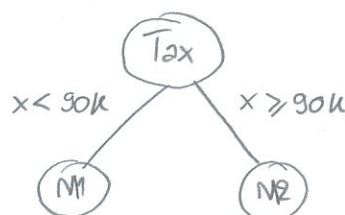
Cheat → No = 1 Cheat → No = 2

$$GINI(M1) = 1 - \left[\left(\frac{0}{1} \right)^2 + \left(\frac{1}{1} \right)^2 \right] = 0$$

$$GINI(M2) = 1 - \left[\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right] = 0.48$$

$$GINI_{split} = \left(\frac{1}{6} \right) (0) + \left(\frac{5}{6} \right) (0.48) = 0.4$$

3.3



Cheat → Yes = 1 Cheat → Yes = 2

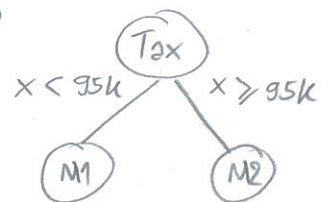
Cheat → No = 1 Cheat → No = 2

$$GINI(M1) = 1 - \left[\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right] = 0.5$$

$$GINI(M2) = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

$$GINI_{split} = \left(\frac{2}{6} \right) (0.5) + \left(\frac{4}{6} \right) (0.5) = 0.5$$

3.4



Cheat → Yes = 2 Cheat → Yes = 1

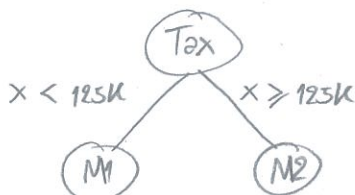
Cheat → No = 1 Cheat → No = 2

$$GINI(M1) = 1 - \left[\left(\frac{2}{3} \right)^2 + \left(\frac{1}{3} \right)^2 \right] = 0.44$$

$$GINI(M2) = 1 - \left[\left(\frac{1}{3} \right)^2 + \left(\frac{2}{3} \right)^2 \right] = 0.44$$

$$GINI_{split} = \left(\frac{3}{6} \right) (0.44) + \left(\frac{3}{6} \right) (0.44) = 0.44$$

3.5



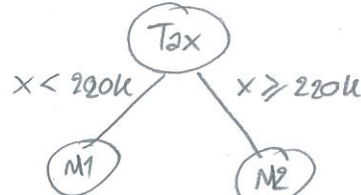
Cheat → Yes = 3 Cheat → Yes = 0

Cheat → No = 1 Cheat → No = 2

$$GINI(M1) = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 0.375$$

$$GINI(M2) = 1 - \left[\left(\frac{0}{2} \right)^2 + \left(\frac{2}{2} \right)^2 \right] = 0$$

$$GINI_{split} = \left(\frac{4}{6} \right) (0.375) + \left(\frac{2}{6} \right) (0) = 0.25$$



Cheat → Yes = 3 Cheat → Yes = 0

Cheat → No = 2 Cheat → No = 1

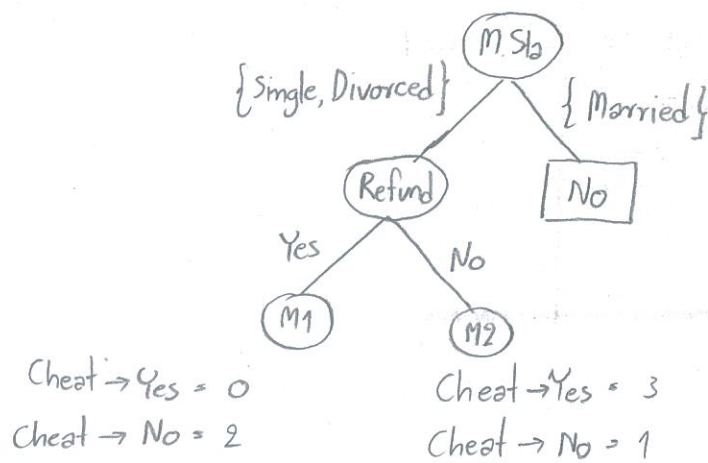
$$GINI(M1) = 1 - \left[\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right] = 0.48$$

$$GINI(M2) = 1 - \left[\left(\frac{0}{1} \right)^2 + \left(\frac{1}{1} \right)^2 \right] = 0$$

$$GINI_{split} = \left(\frac{5}{6} \right) (0.48) + \left(\frac{1}{6} \right) (0) = 0.4$$

Minimum { 0.25, 0.5, 0.5, 0.4, 0.5, 0.44, 0.25, 0.4 } = 0.25

Since the minimum $GINI_{split}$ of 0.25 is obtained from test condition 1 and 3.5, thus we can use either test condition 1 or 3.5 as a test condition for node N1. In this case, I choose test condition 1.



Step 5 : Check whether we should stop splitting the records at node M1.

\therefore Since all the records at node M1 satisfy the first main stopping criterion, then we stop splitting and declare this node as a leaf node of value No.

Step 6 : Check whether we should stop splitting the records at node M2.

\therefore Since all the records at node M2 do not satisfy the two main stopping criteria, then we need to split these records further.

Repeat step 4, 5, and 6 until we cannot split further.

