# Decision Tree Analysis Based on Entropy Values
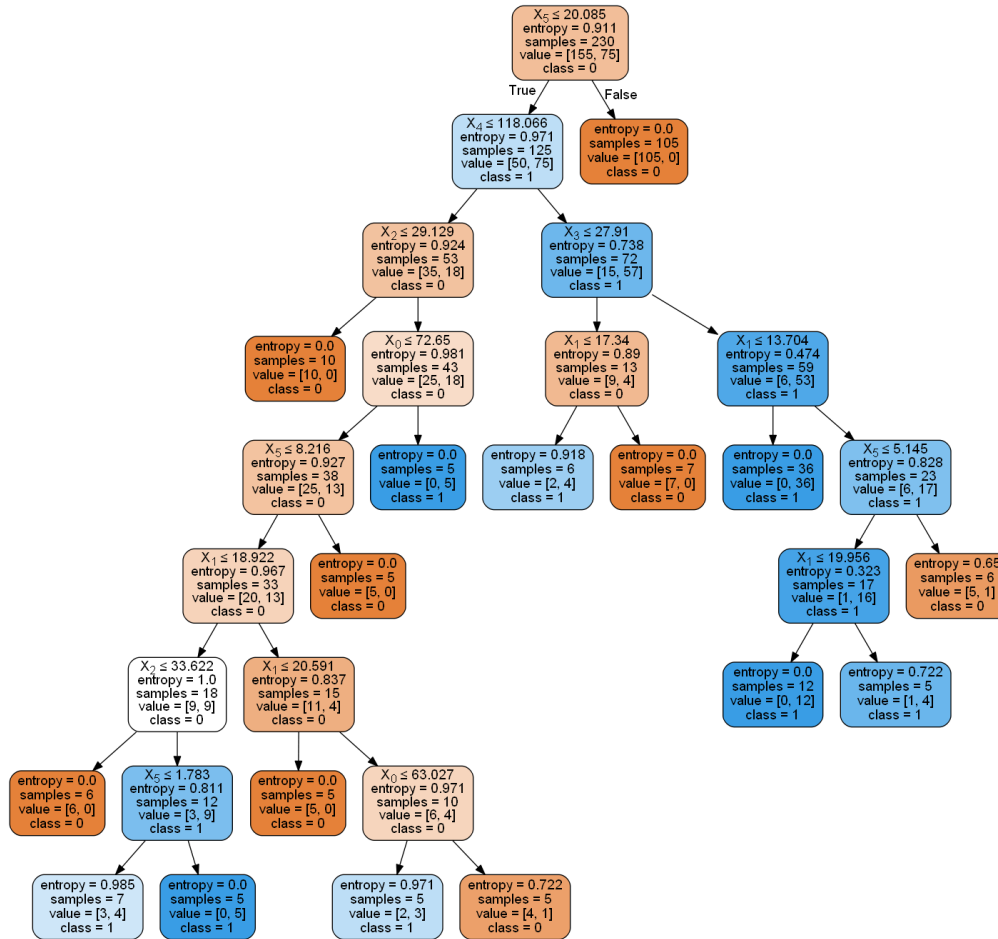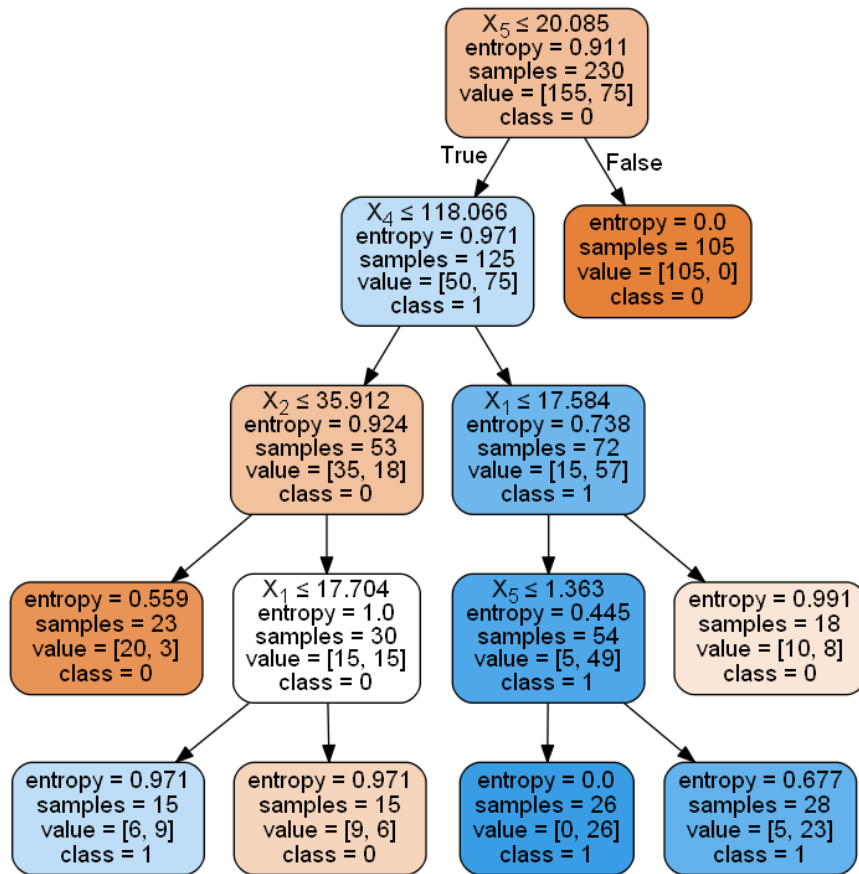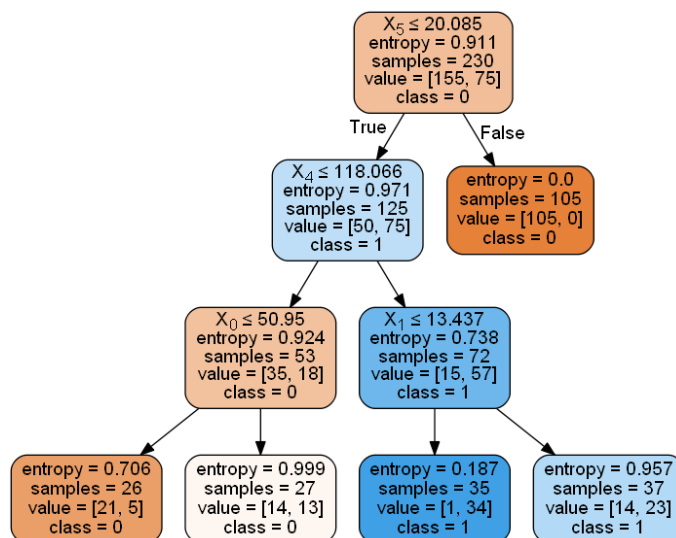
**Minimum 5 sample leaf nodes**
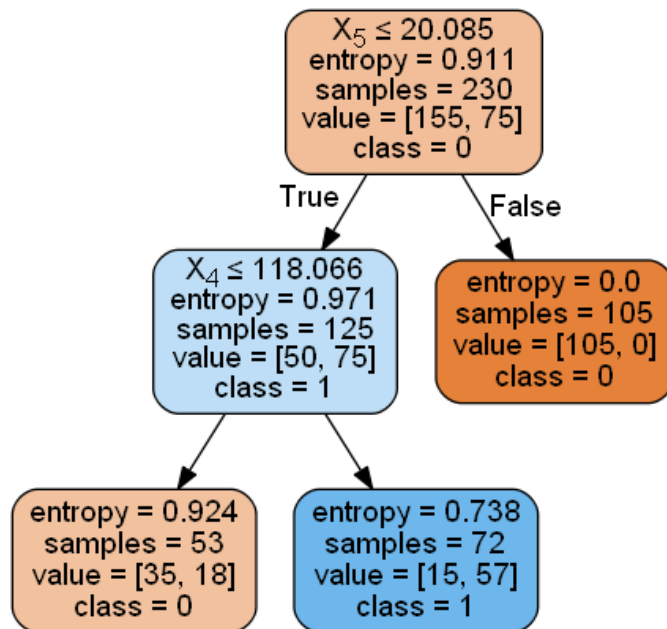
## 15 samples leaf node

**Root:** $X_5 \le 20.085$
entropy = 0.911
samples = 230
value = [155, 75]
class = 0

True → **$X_4 \le 118.066$**
entropy = 0.971
samples = 125
value = [50, 75]
class = 1

False → entropy = 0.0
samples = 105
value = [105, 0]
class = 0

**$X_2 \le 35.912$**
entropy = 0.924
samples = 53
value = [35, 18]
class = 0

**$X_1 \le 17.584$**
entropy = 0.738
samples = 72
value = [15, 57]
class = 1

entropy = 0.559
samples = 23
value = [20, 3]
class = 0

**$X_1 \le 17.704$**
entropy = 1.0
samples = 30
value = [15, 15]
class = 0

**$X_5 \le 1.363$**
entropy = 0.445
samples = 54
value = [5, 49]
class = 1

entropy = 0.991
samples = 18
value = [10, 8]
class = 0

entropy = 0.971
samples = 15
value = [6, 9]
class = 1

entropy = 0.971
samples = 15
value = [9, 6]
class = 0

entropy = 0.0
samples = 26
value = [0, 26]
class = 1

entropy = 0.677
samples = 28
value = [5, 23]
class = 1

## 25 sample leaf node

**Root:** $X_5 \le 20.085$
entropy = 0.911
samples = 230
value = [155, 75]
class = 0

True → **$X_4 \le 118.066$**
entropy = 0.971
samples = 125
value = [50, 75]
class = 1

False → entropy = 0.0
samples = 105
value = [105, 0]
class = 0

**$X_0 \le 50.95$**
entropy = 0.924
samples = 53
value = [35, 18]
class = 0

**$X_1 \le 13.437$**
entropy = 0.738
samples = 72
value = [15, 57]
class = 1

entropy = 0.706
samples = 26
value = [21, 5]
class = 0

entropy = 0.999
samples = 27
value = [14, 13]
class = 0

entropy = 0.187
samples = 35
value = [1, 34]
class = 1

entropy = 0.957
samples = 37
value = [14, 23]
class = 1
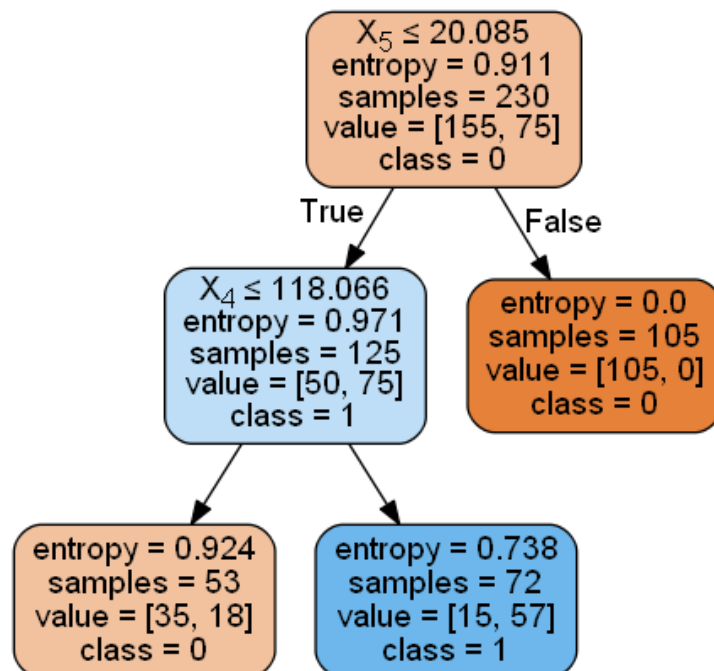
**Min 40 sample leaf node:**



**Min 50 sample leaf node**

**Conclusion:**

Among the 5 decision trees, I would pick the decision tree with the min 15 samples in the leaf node. As we increase min sample value in leaf node the decision tree shallower, However, if we check the impurity value entropy values are higher in the shallow trees. This shows that we can't achieve our goal maximizing the information gain. Although the tree with min 5 sample in leaf node gives the good information gain, It is not very shallow and simple. I don't see huge information gain difference between 5 sample leaf and 15 sample leaf. As considering Occam's Razor, simpler model rule. I chose the tree with min 15 samples in its leaf nodes.

**Accuracy and Precision Values**

**Minimum 5 sample leaf nodes**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.88      | 0.82   | 0.85     | 55      |
| 1            | 0.66      | 0.76   | 0.70     | 25      |
| accuracy     |           |        | 0.80     | 80      |
| macro avg    | 0.77      | 0.79   | 0.78     | 80      |
| weighted avg | 0.81      | 0.80   | 0.80     | 80      |

**Minimum 15 sample leaf nodes**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.87      | 0.84   | 0.85     | 55      |
| 1            | 0.67      | 0.72   | 0.69     | 25      |
| accuracy     |           |        | 0.80     | 80      |
| macro avg    | 0.77      | 0.78   | 0.77     | 80      |
| weighted avg | 0.81      | 0.80   | 0.80     | 80      |

**Minimum 25 sample leaf nodes**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.86      | 0.78   | 0.82     | 55      |
| 1            | 0.60      | 0.72   | 0.65     | 25      |
| accuracy     |           |        | 0.76     | 80      |
| macro avg    | 0.73      | 0.75   | 0.74     | 80      |
| weighted avg | 0.78      | 0.76   | 0.77     | 80      |

**Minimum 40 sample leaf nodes**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.86      | 0.78   | 0.82     | 55      |
| 1            | 0.60      | 0.72   | 0.65     | 25      |
|              |           |        |          |         |
| accuracy     |           |        | 0.76     | 80      |
| macro avg    | 0.73      | 0.75   | 0.74     | 80      |
| weighted avg | 0.78      | 0.76   | 0.77     | 80      |

**Minimum 50 sample leaf nodes**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.86      | 0.78   | 0.82     | 55      |
| 1            | 0.60      | 0.72   | 0.65     | 25      |
|              |           |        |          |         |
| accuracy     |           |        | 0.76     | 80      |
| macro avg    | 0.73      | 0.75   | 0.74     | 80      |
| weighted avg | 0.78      | 0.76   | 0.77     | 80      |

5 and 15 sample leaf node models have the same accuracy rate 80% while others show the 76 % of accuracy. Moreover,  5 and 15 sample leaf node models are higher precision and recall values that give us better exactness and relevancy.


**Repeating Analysis with Having Three Different Class Labels**

**Min 5 sample leaf node**



Decision tree (Min 5 sample leaf node):

- $X_5 \le 15.153$, entropy = 1.497, samples = 230, value = [45, 75, 110], class = 2
  - True → $X_4 \le 118.066$, entropy = 1.1, samples = 122, value = [45, 74, 3], class = 1
  - False → $X_5 \le 21.258$, entropy = 0.076, samples = 108, value = [0, 1, 107], class = 2
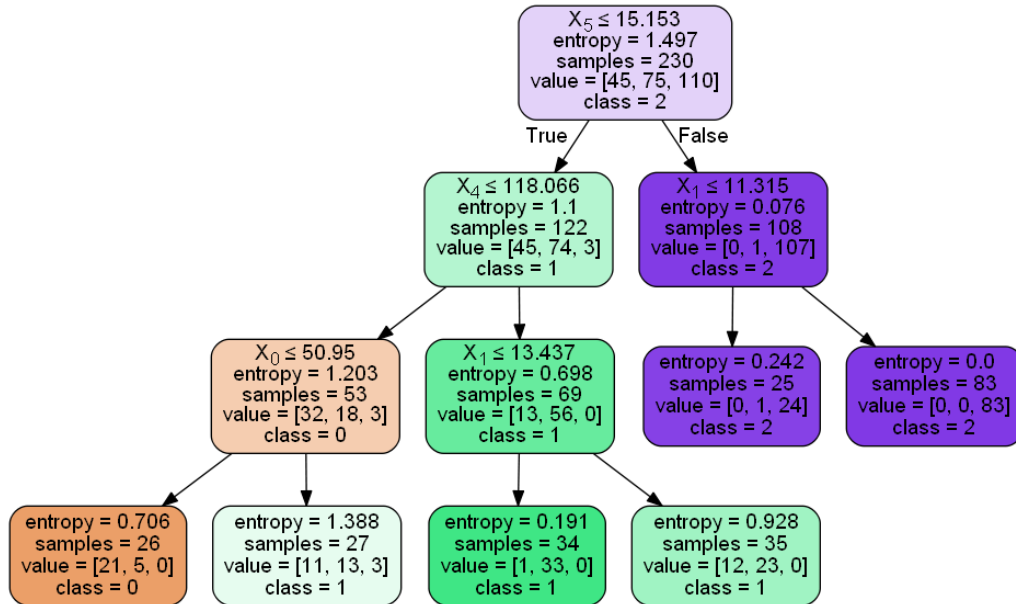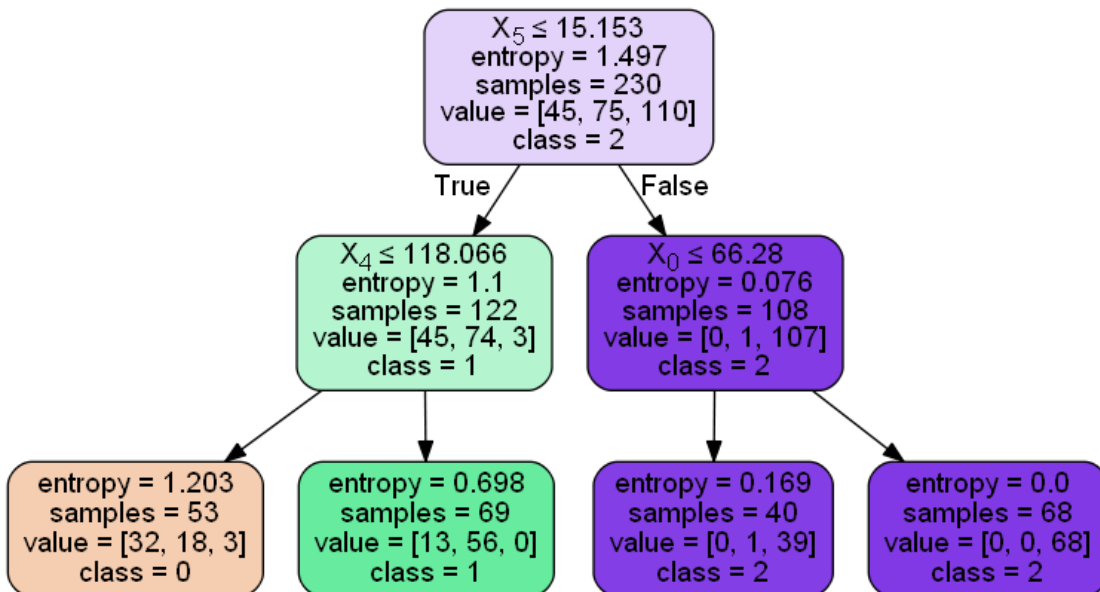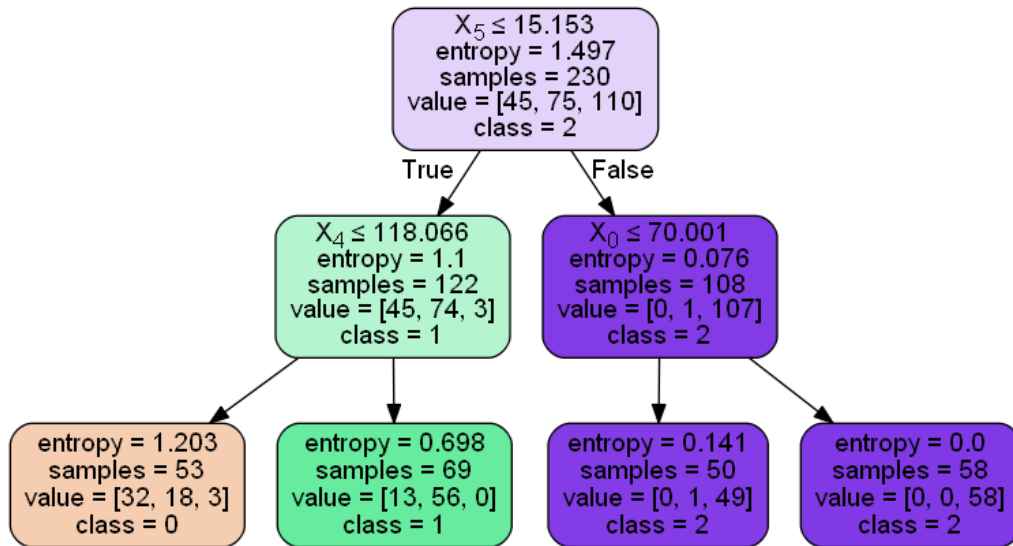    - entropy = 0.722, samples = 5, value = [0, 1, 4], class = 2
    - entropy = 0.0, samples = 103, value = [0, 0, 103], class = 2
  - $X_3 \le 47.149$, entropy = 1.203, samples = 53, value = [32, 18, 3], class = 0
  - $X_3 \le 27.91$, entropy = 0.698, samples = 69, value = [13, 56, 0], class = 1
    - $X_4 \le 113.628$, entropy = 0.958, samples = 44, value = [32, 11, 1], class = 0
    - entropy = 0.764, samples = 9, value = [0, 7, 2], class = 1
    - $X_1 \le 17.34$, entropy = 0.89, samples = 13, value = [9, 4, 0], class = 0
    - $X_1 \le 14.176$, entropy = 0.371, samples = 56, value = [4, 52, 0], class = 1
      - $X_1 \le 22.386$, entropy = 0.297, samples = 19, value = [18, 1, 0], class = 0
      - $X_2 \le 29.129$, entropy = 1.183, samples = 25, value = [14, 10, 1], class = 0
      - entropy = 0.918, samples = 6, value = [2, 4, 0], class = 1
      - entropy = 0.0, samples = 7, value = [7, 0, 0], class = 0
      - entropy = 0.0, samples = 37, value = [0, 37, 0], class = 1
      - $X_5 \le 4.914$, entropy = 0.742, samples = 19, value = [4, 15, 0], class = 1
        - entropy = 0.0, samples = 14, value = [14, 0, 0], class = 0
        - entropy = 0.722, samples = 5, value = [4, 1, 0], class = 0
        - entropy = 0.0, samples = 5, value = [5, 0, 0], class = 0
        - $X_4 \le 115.84$, entropy = 1.234, samples = 20, value = [9, 10, 1], class = 1
        - $X_4 \le 121.328$, entropy = 0.371, samples = 14, value = [1, 13, 0], class = 1
        - entropy = 0.971, samples = 5, value = [3, 2, 0], class = 0
          - entropy = 0.592, samples = 7, value = [1, 6, 0], class = 1
          - $X_0 \le 48.118$, entropy = 1.239, samples = 13, value = [8, 4, 1], class = 0
          - entropy = 0.722, samples = 5, value = [1, 4, 0], class = 1
          - entropy = 0.0, samples = 9, value = [0, 9, 0], class = 1
            - entropy = 1.0, samples = 8, value = [4, 4, 0], class = 0
            - entropy = 0.722, samples = 5, value = [4, 0, 1], class = 0

**Min 15 samples leaf node**

Decision tree (Min 15 samples leaf node):

- $X_5 \le 15.153$, entropy = 1.497, samples = 230, value = [45, 75, 110], class = 2
  - True → $X_4 \le 118.066$, entropy = 1.1, samples = 122, value = [45, 74, 3], class = 1
  - False → $X_5 \le 25.656$, entropy = 0.076, samples = 108, value = [0, 1, 107], class = 2
    - $X_2 \le 48.139$, entropy = 1.203, samples = 53, value = [32, 18, 3], class = 0
    - $X_3 \le 30.001$, entropy = 0.698, samples = 69, value = [13, 56, 0], class = 1
    - entropy = 0.353, samples = 15, value = [0, 1, 14], class = 2
    - entropy = 0.0, samples = 93, value = [0, 0, 93], class = 2
      - $X_4 \le 113.523$, entropy = 0.831, samples = 38, value = [28, 10, 0], class = 0
      - entropy = 1.457, samples = 15, value = [4, 8, 3], class = 1
      - entropy = 0.991, samples = 18, value = [10, 8, 0], class = 0
      - $X_4 \le 121.725$, entropy = 0.323, samples = 51, value = [3, 48, 0], class = 1
        - entropy = 0.337, samples = 16, value = [15, 1, 0], class = 0
        - entropy = 0.976, samples = 22, value = [13, 9, 0], class = 0
        - entropy = 0.672, samples = 17, value = [3, 14, 0], class = 1
        - entropy = 0.0, samples = 34, value = [0, 34, 0], class = 1

**Min 25 samples leaf node**

$X_5 \leq 15.153$
entropy = 1.497
samples = 230
value = [45, 75, 110]
class = 2

True / False

$X_4 \leq 118.066$
entropy = 1.1
samples = 122
value = [45, 74, 3]
class = 1

$X_1 \leq 11.315$
entropy = 0.076
samples = 108
value = [0, 1, 107]
class = 2

$X_0 \leq 50.95$
entropy = 1.203
samples = 53
value = [32, 18, 3]
class = 0

$X_1 \leq 13.437$
entropy = 0.698
samples = 69
value = [13, 56, 0]
class = 1

entropy = 0.242
samples = 25
value = [0, 1, 24]
class = 2

entropy = 0.0
samples = 83
value = [0, 0, 83]
class = 2

entropy = 0.706
samples = 26
value = [21, 5, 0]
class = 0

entropy = 1.388
samples = 27
value = [11, 13, 3]
class = 1

entropy = 0.191
samples = 34
value = [1, 33, 0]
class = 1

entropy = 0.928
samples = 35
value = [12, 23, 0]
class = 1

**Min 40 samples leaf node**

$X_5 \leq 15.153$
entropy = 1.497
samples = 230
value = [45, 75, 110]
class = 2

True / False

$X_4 \leq 118.066$
entropy = 1.1
samples = 122
value = [45, 74, 3]
class = 1

$X_0 \leq 66.28$
entropy = 0.076
samples = 108
value = [0, 1, 107]
class = 2

entropy = 1.203
samples = 53
value = [32, 18, 3]
class = 0

entropy = 0.698
samples = 69
value = [13, 56, 0]
class = 1

entropy = 0.169
samples = 40
value = [0, 1, 39]
class = 2

entropy = 0.0
samples = 68
value = [0, 0, 68]
class = 2

**Min 50 samples leaf node**



**Conclusion:**

Min 5 sample decision tree is shallower in this decision tree. I think, adding one more class helps for classifying better in this data set. On the other hand, in each every decision trees' leaf node, entrophy is very high for the class 0. So we don't have information gain for class 0. In every decision tree, class 2 is easily splitted. Also, compared to the decision tree in the previous example, I would still select the decision tree with min 15 sample leaf node  because of the precision recall  accuracy and weighted average values. I was leaning towards to select 25 samples one however its weighted average value and accuracy is very low.

**Accuracy and Precision Values for 3 different class labels scenario**

**Min 5 samples leaf node**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.50 | 0.47 | 0.48 | 15 |
| 1 | 0.68 | 0.68 | 0.68 | 25 |
| 2 | 0.93 | 0.95 | 0.94 | 40 |
| accuracy |  |  | 0.78 | 80 |
| macro avg | 0.70 | 0.70 | 0.70 | 80 |
| weighted avg | 0.77 | 0.78 | 0.77 | 80 |

**Min 15 samples leaf node**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.59 | 0.67 | 0.62 | 15 |

```
     1     0.73    0.64    0.68     25
     2     0.93    0.95    0.94     40

  accuracy                  0.80     80
 macro avg     0.75    0.75    0.75     80
weighted avg     0.80    0.80    0.80     80
```

**Min 25 samples leaf node**

```
         precision   recall  f1-score   support

     0     0.08    0.27    0.12     15
     1     0.60    0.72    0.65     25
     2     0.00    0.00    0.00     40

  accuracy                  0.28     80
 macro avg     0.23    0.33    0.26     80
weighted avg     0.20    0.28    0.23     80
```

**Min 40 samples leaf node**

```
         precision   recall  f1-score   support

     0     0.33    0.20    0.25     15
     1     0.60    0.72    0.65     25
     2     0.93    0.95    0.94     40

  accuracy                  0.74     80
 macro avg     0.62    0.62    0.61     80
weighted avg     0.71    0.74    0.72     80
```

**Min 50 samples leaf node**

```
         precision   recall  f1-score   support

     0     0.33    0.20    0.25     15
     1     0.60    0.72    0.65     25
     2     0.93    0.95    0.94     40

  accuracy                  0.74     80
 macro avg     0.62    0.62    0.61     80
weighted avg     0.71    0.74    0.72     80
```

**Eliminating one of the variable that has high correlation with another independent variable and repeating decision tree analysis**
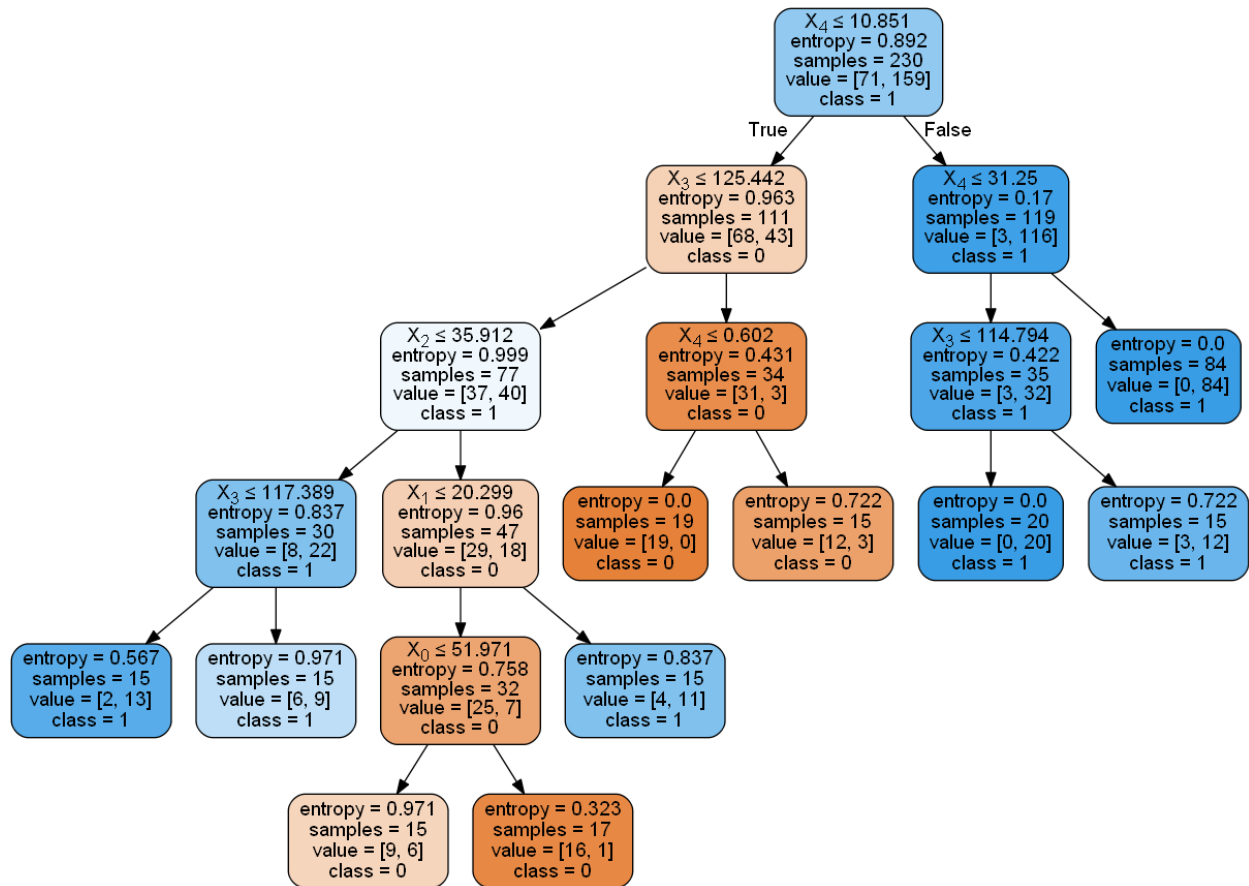


corr - DataFrame

| Index | pelvic_incidence | pelvic_tilt numeric | nbar_lordosis_ang | sacral_slope | pelvic_radius | ree_spondylolisth | diagnosed |
|---|---|---|---|---|---|---|---|
| pelvic_incide... | 1.000 | 0.629 | 0.717 | 0.815 | -0.247 | 0.639 | 0.353 |
| pelvic_tilt numeric | 0.629 | 1.000 | 0.433 | 0.062 | 0.033 | 0.398 | 0.326 |
| lumbar_lordos... | 0.717 | 0.433 | 1.000 | 0.598 | -0.080 | 0.534 | 0.312 |
| sacral_slope | 0.815 | 0.062 | 0.598 | 1.000 | -0.342 | 0.524 | 0.211 |
| pelvic_radius | -0.247 | 0.033 | -0.080 | -0.342 | 1.000 | -0.026 | -0.310 |
| degree_spondy... | 0.639 | 0.398 | 0.534 | 0.524 | -0.026 | 1.000 | 0.444 |
| diagnosed | 0.353 | 0.326 | 0.312 | 0.211 | -0.310 | 0.444 | 1.000 |

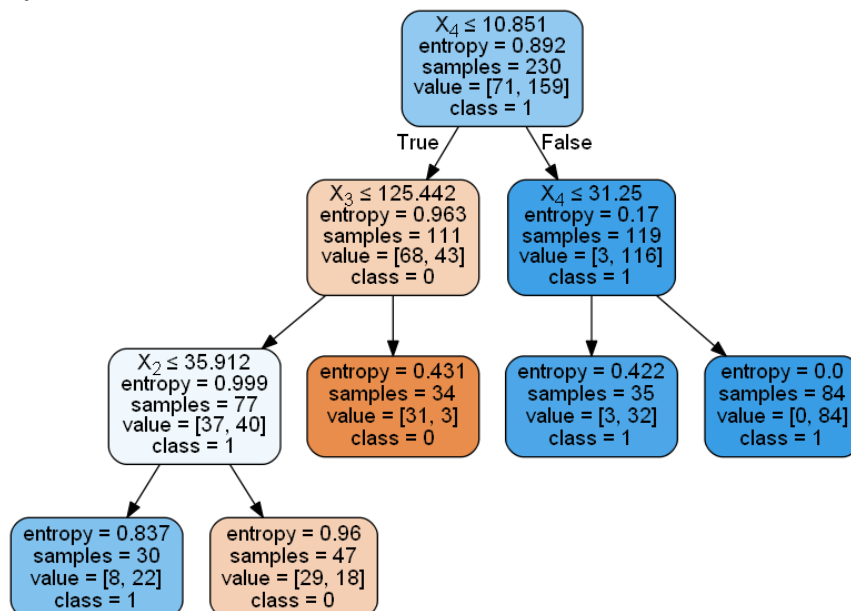I eliminated the variable sacral_slope since it has high correaltion with pelvis_incidence.
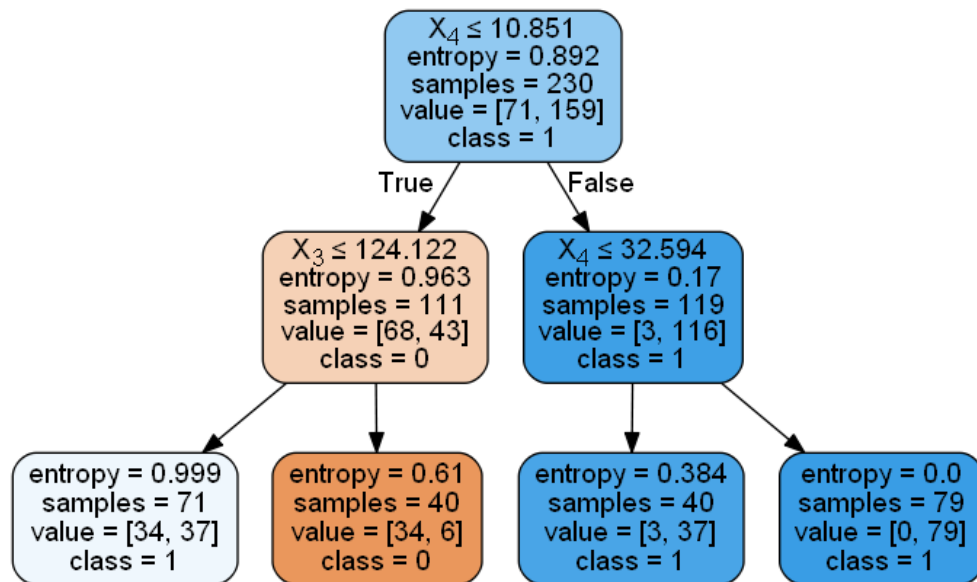
**Min 5 samples leaf node**

**Min 15 samples leaf node**

$X_4 \leq 10.851$
entropy = 0.892
samples = 230
value = [71, 159]
class = 1

True / False

$X_3 \leq 125.442$
entropy = 0.963
samples = 111
value = [68, 43]
class = 0

$X_4 \leq 31.25$
entropy = 0.17
samples = 119
value = [3, 116]
class = 1

$X_2 \leq 35.912$
entropy = 0.999
samples = 77
value = [37, 40]
class = 1

$X_4 \leq 0.602$
entropy = 0.431
samples = 34
value = [31, 3]
class = 0

$X_3 \leq 114.794$
entropy = 0.422
samples = 35
value = [3, 32]
class = 1

entropy = 0.0
samples = 84
value = [0, 84]
class = 1

$X_3 \leq 117.389$
entropy = 0.837
samples = 30
value = [8, 22]
class = 1

$X_1 \leq 20.299$
entropy = 0.96
samples = 47
value = [29, 18]
class = 0

entropy = 0.0
samples = 19
value = [19, 0]
class = 0

entropy = 0.722
samples = 15
value = [12, 3]
class = 0

entropy = 0.0
samples = 20
value = [0, 20]
class = 1

entropy = 0.722
samples = 15
value = [3, 12]
class = 1

entropy = 0.567
samples = 15
value = [2, 13]
class = 1

entropy = 0.971
samples = 15
value = [6, 9]
class = 1

$X_0 \leq 51.971$
entropy = 0.758
samples = 32
value = [25, 7]
class = 0

entropy = 0.837
samples = 15
value = [4, 11]
class = 1

entropy = 0.971
samples = 15
value = [9, 6]
class = 0

entropy = 0.323
samples = 17
value = [16, 1]
class = 0

**Min 25 samples leaf node**

$X_4 \leq 10.851$
entropy = 0.892
samples = 230
value = [71, 159]
class = 1

True / False

$X_3 \leq 125.442$
entropy = 0.963
samples = 111
value = [68, 43]
class = 0

$X_4 \leq 31.25$
entropy = 0.17
samples = 119
value = [3, 116]
class = 1

$X_2 \leq 35.912$
entropy = 0.999
samples = 77
value = [37, 40]
class = 1

entropy = 0.431
samples = 34
value = [31, 3]
class = 0

entropy = 0.422
samples = 35
value = [3, 32]
class = 1

entropy = 0.0
samples = 84
value = [0, 84]
class = 1

entropy = 0.837
samples = 30
value = [8, 22]
class = 1

entropy = 0.96
samples = 47
value = [29, 18]
class = 0

## Min 40 samples leaf node

```
                    X₄ ≤ 10.851
                   entropy = 0.892
                   samples = 230
                  value = [71, 159]
                      class = 1
           True  /                  \  False
                /                      \
      X₃ ≤ 124.122                   X₄ ≤ 32.594
     entropy = 0.963               entropy = 0.17
     samples = 111                 samples = 119
     value = [68, 43]              value = [3, 116]
        class = 0                     class = 1
      /           \                 /            \
entropy=0.999  entropy=0.61   entropy=0.384  entropy=0.0
samples=71     samples=40     samples=40     samples=79
value=[34,37]  value=[34,6]   value=[3,37]   value=[0,79]
class=1        class=0        class=1        class=1
```

## Min 50 samples leaf node

```
                    X₄ ≤ 10.851
                   entropy = 0.892
                   samples = 230
                  value = [71, 159]
                      class = 1
           True  /                  \  False
                /                      \
      X₃ ≤ 119.354                   X₄ ≤ 37.795
     entropy = 0.963               entropy = 0.17
     samples = 111                 samples = 119
     value = [68, 43]              value = [3, 116]
        class = 0                     class = 1
      /           \                 /            \
entropy=0.983  entropy=0.761  entropy=0.327  entropy=0.0
samples=52     samples=59     samples=50     samples=69
value=[22,30]  value=[46,13]  value=[3,47]   value=[0,69]
class=1        class=0        class=1        class=1
```

## Accuracy and Precision Values

### Min 5 samples leaf node

|   |      |      |      |    |
|---|------|------|------|----|
| 0 | 0.74 | 0.69 | 0.71 | 29 |
| 1 | 0.83 | 0.86 | 0.85 | 51 |

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| accuracy     |           |        | 0.80     | 80      |
| macro avg    | 0.79      | 0.78   | 0.78     | 80      |
| weighted avg | 0.80      | 0.80   | 0.80     | 80      |

**Min 15 samples leaf node**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.75      | 0.62   | 0.68     | 29      |
| 1            | 0.80      | 0.88   | 0.84     | 51      |
| accuracy     |           |        | 0.79     | 80      |
| macro avg    | 0.78      | 0.75   | 0.76     | 80      |
| weighted avg | 0.78      | 0.79   | 0.78     | 80      |

**Min 25 samples leaf node**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.71      | 0.69   | 0.70     | 29      |
| 1            | 0.83      | 0.84   | 0.83     | 51      |
| accuracy     |           |        | 0.79     | 80      |
| macro avg    | 0.77      | 0.77   | 0.77     | 80      |
| weighted avg | 0.79      | 0.79   | 0.79     | 80      |

**Min 50 samples leaf node**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.73      | 0.38   | 0.50     | 29      |
| 1            | 0.72      | 0.92   | 0.81     | 51      |
| accuracy     |           |        | 0.73     | 80      |
| macro avg    | 0.73      | 0.65   | 0.66     | 80      |
| weighted avg | 0.73      | 0.72   | 0.70     | 80      |

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.69      | 0.62   | 0.65     | 29      |
| 1            | 0.80      | 0.84   | 0.82     | 51      |
| accuracy     |           |        | 0.76     | 80      |
| macro avg    | 0.74      | 0.73   | 0.74     | 80      |
| weighted avg | 0.76      | 0.76   | 0.76     | 80      |

**Conclusion :**
When I left the variable sacral_slope, accuracy increased for decision tree with 25 min sample leaf and above others. Also, the decision tree depth for 5 min sample became little shallower. Comparatively, entropy values are in leaves are smaller than before we eliminated the correlated value. Since our sample size is not large, the change in the accuracy is not very big but we can still observe it.