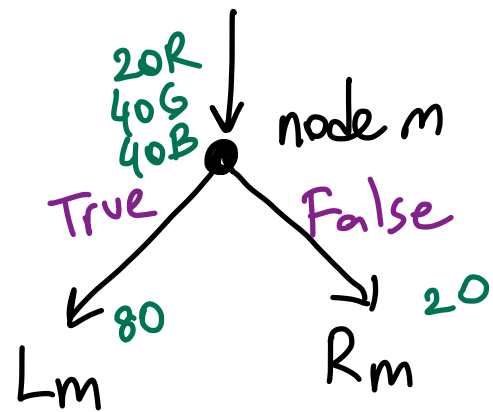


Univariate Decision Trees

$$L_m = \{x \mid x_j > \widetilde{w_{m0}}\}$$

threshold
feature index

$$R_m = \{x \mid x_j \leq w_{m0}\}$$



$$P_{m1} = \frac{20}{100}$$

$$P_{m2} = \frac{40}{100}$$

$$P_{m3} = \frac{40}{100}$$

$N_m = 100 \leftarrow N_m = \# \text{ of data points that reach node } m$

$N_{m,c} = \# \text{ of data points that reach node } m \text{ from class } \#c$
 $k=3$ $N_{m,1}=20$ $N_{m,2}=40$ $N_{m,3}=40$

$N_{m,s} = \# \text{ of data points that reach node } m \text{ and takes split } \#s$
 $S=2$ $N_{m,1}=80$ $N_{m,2}=20$

impurity of a node

$$P_{mc} = \hat{P}(y=c \mid \underline{x_m}) = \frac{N_{m,c}}{N_m}$$

$$I_m = - \sum_{c=1}^K P_{mc} \cdot \log_2(P_{mc})$$

of classes

$$\frac{80}{100} I_m(L_m) + \frac{20}{100} I_m(R_m)$$

$$I_m' = \sum_{s=1}^S \left[\underbrace{\frac{N_{m,s}}{N_m}}_{\text{weight of child node}} \left[\underbrace{- \sum_{c=1}^K P_{msc} \log_2(P_{msc})}_{\text{impurity of child node}} \right] \right]$$

impurity of a split

Entropy:

$$-p \log_2(p) - (1-p) \log_2(1-p)$$

$$0 \log_2(0) \equiv 0$$

p = ratio of positive data points
 $1-p$ = ratio of negative data points

$$\frac{N_+}{N_+ + N_-}$$
$$\frac{N_-}{N_+ + N_-}$$

Gini Index:

$$2 \cdot p \cdot (1-p)$$

Gini Index

1/2



all data points are negative

all data points are positive

Misclassification Error:

$$1 - \max(p, 1-p)$$

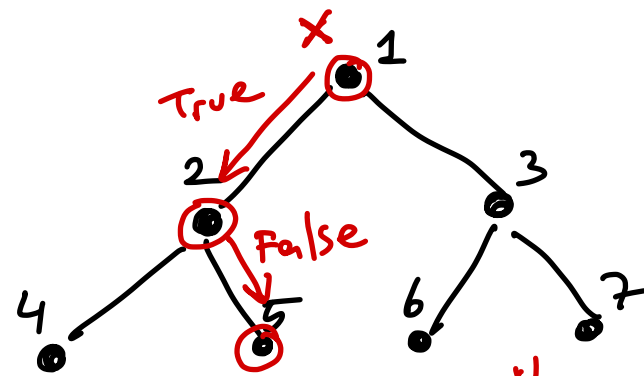
60% 40%

$$\text{OR } \min(p, 1-p)$$

multiclass classification $\Rightarrow 1 - \max(p_1, p_2, \dots, p_k)$
classification accuracy when majority label is used

Regression Trees:

$$b_m(x) = \begin{cases} 1 & \text{if } x \in \mathcal{X}_m \text{ (x reaches node m)} \\ 0 & \text{otherwise} \end{cases}$$



$$\begin{aligned} b_1(x) &= 1 & b_2(x) &= 1 & b_3(x) &= 0 \\ b_4(x) &= 0 & b_5(x) &= 1 & b_6(x) &= 0 \\ b_7(x) &= 0 \end{aligned}$$

$$N_m = \sum_{i=1}^N b_m(x_i)$$

error of a node.

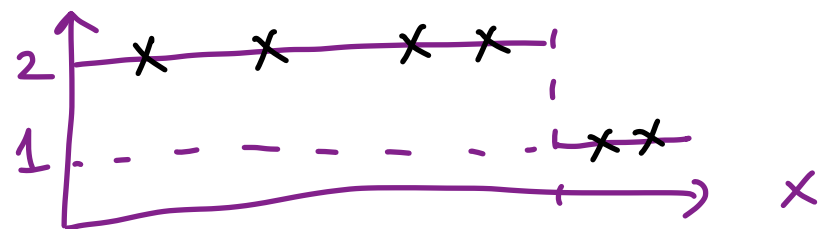
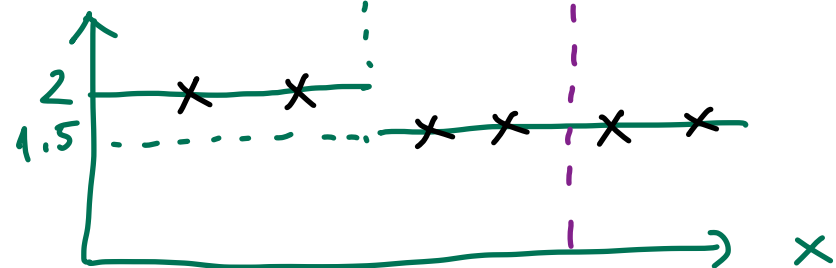
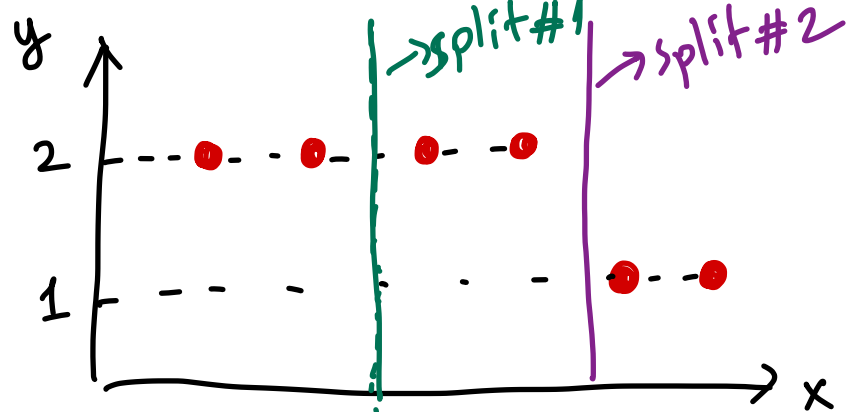
$$E_m = \frac{1}{N_m} \sum_{i=1}^N \left[(y_i - g_m)^2 b_m(x_i) \right]$$

→ predicted value at node m.
→ # of data points that reach to node m.

$$g_m = \frac{\sum_{i=1}^N [y_i b_m(x_i)]}{\sum_{i=1}^N b_m(x_i)} \quad \left. \vphantom{\sum_{i=1}^N} \right\} \text{average response (sample mean)}$$

error of a split

$$E_m = \frac{1}{N_m} \sum_{s=1}^S \sum_{i=1}^N \left[(y_i - g_{ms})^2 b_{ms}(x_i) \right]$$



$$E(S_1) = \frac{1}{6} \left[(2-2)^2 + (2-2)^2 + (2-1.5)^2 + (2-1.5)^2 + (1-1.5)^2 + (1-1.5)^2 \right] = 1/6$$

$$E(S_2) = \frac{1}{6} \left[(2-2)^2 + (2-2)^2 + (2-2)^2 + (2-2)^2 + (1-1)^2 + (1-1)^2 \right] = 0$$

$$W = [w_1, w_2, \dots, w_D]$$

$$\begin{bmatrix} 0 & 0 & \dots & 1 & \dots & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_j \\ \vdots \\ x_D - 1 \\ x_D \end{bmatrix} - w_{m0} > 0$$

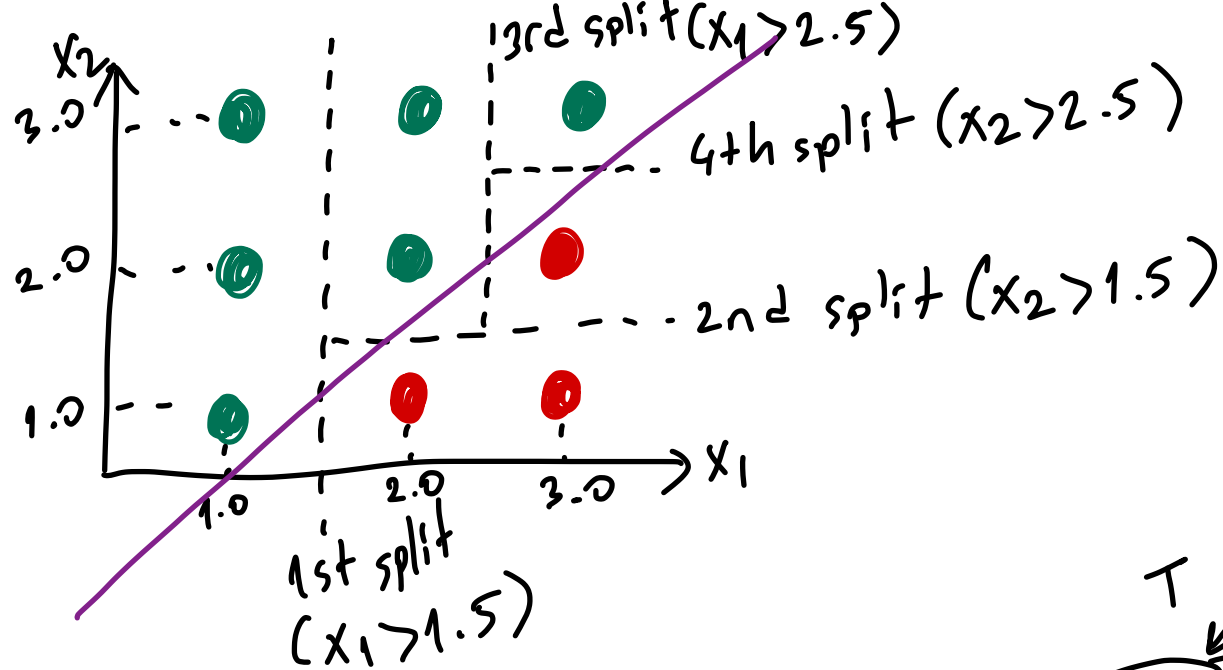
Multivariate Decision Trees

$$f_m(x): x_j > w_{m0}$$

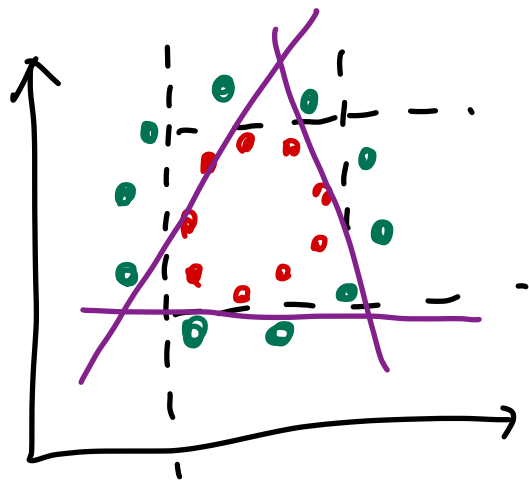
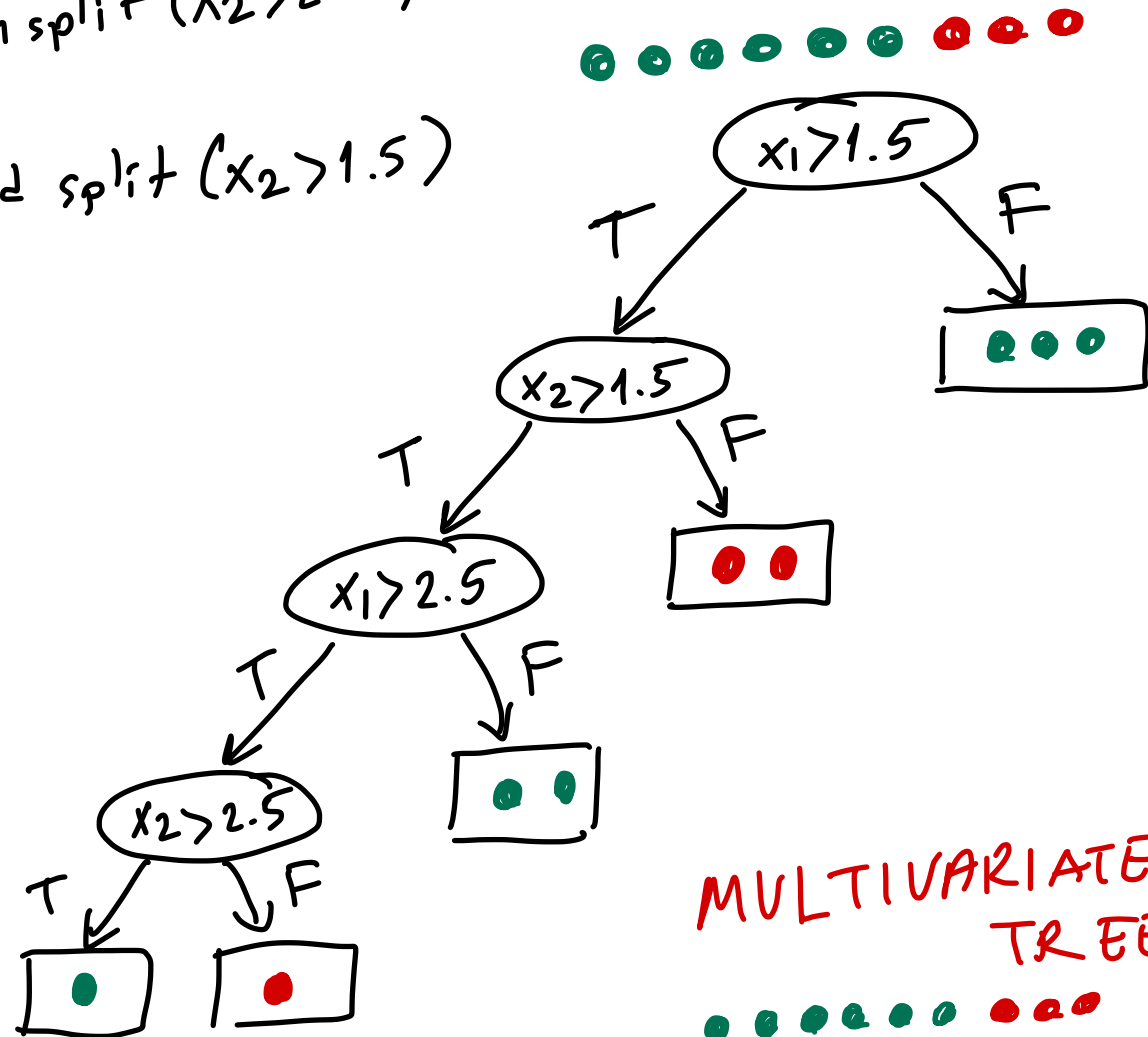
$$f_m(x): \underbrace{w_m^T \cdot x}_{\text{arbitrary vector in } \mathbb{R}^D} + \underbrace{w_{m0}}_{\text{threshold}} > 0$$

$$\Rightarrow x_j - w_{m0} > 0 \Leftarrow \text{univariate split}$$

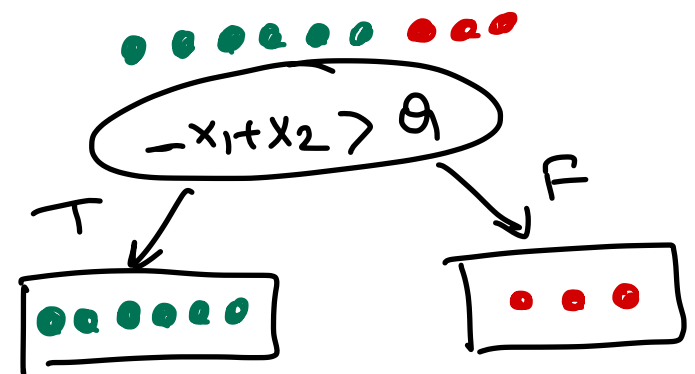
$$\Leftarrow \text{multivariate split.}$$



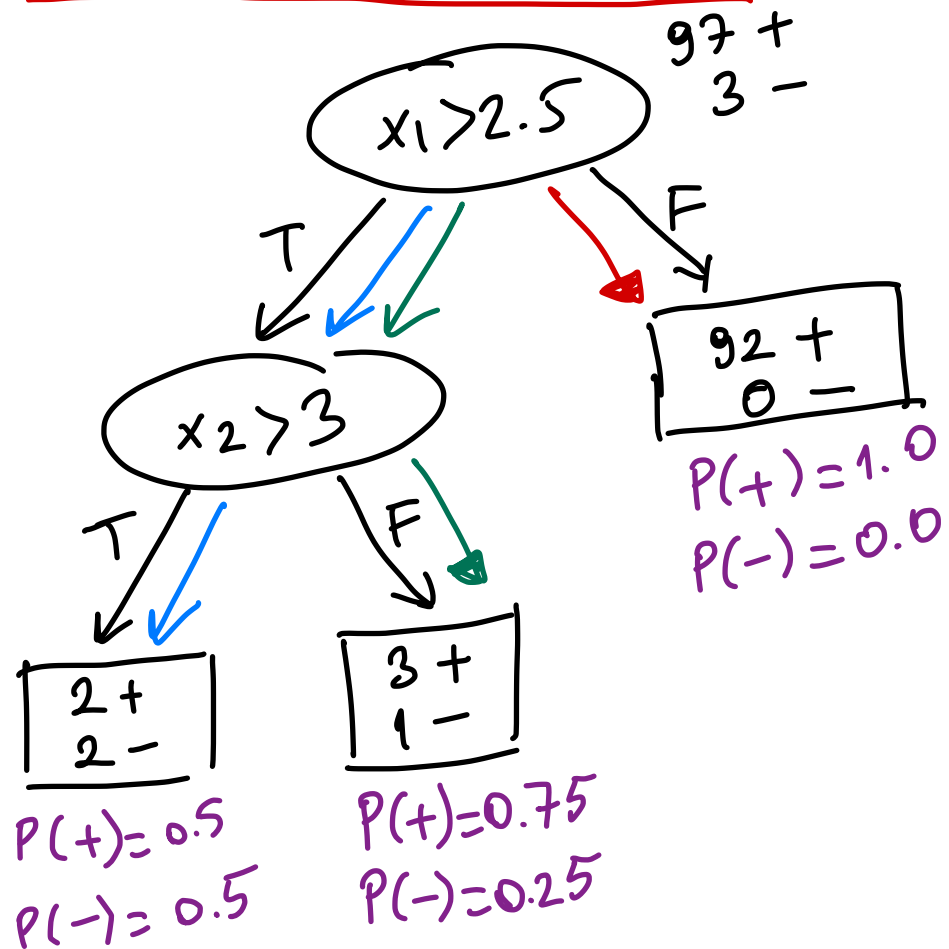
UNIVARIATE TREE



MULTIVARIATE TREE



RULE EXTRACTION

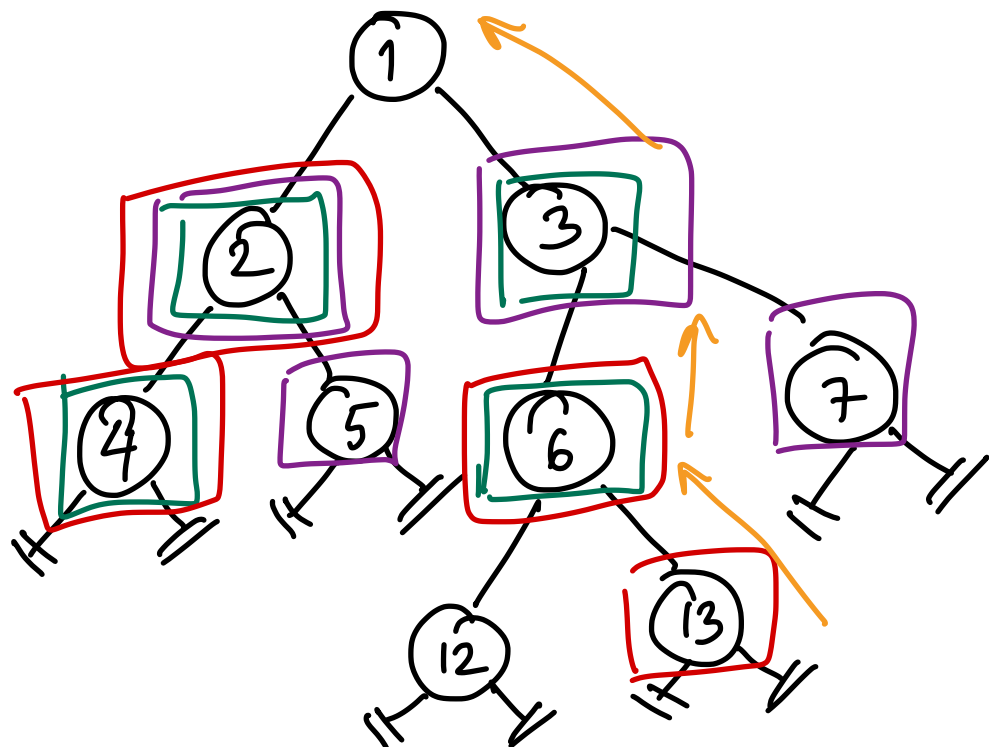


-extract one rule set
for each terminal node.

Path 1: $x_1 \leq 2.5$

Path 2: $x_1 > 2.5$ AND $x_2 > 3$

Path 3: $x_1 > 2.5$ AND $x_2 \leq 3$



$$\text{left child} = 2 * \text{parent}$$

$$6 = 2 * 3$$

$$4 = 2 * 2$$

$$\text{right child} = 2 * \text{parent} + 1$$

$$7 = 2 * 3 + 1$$

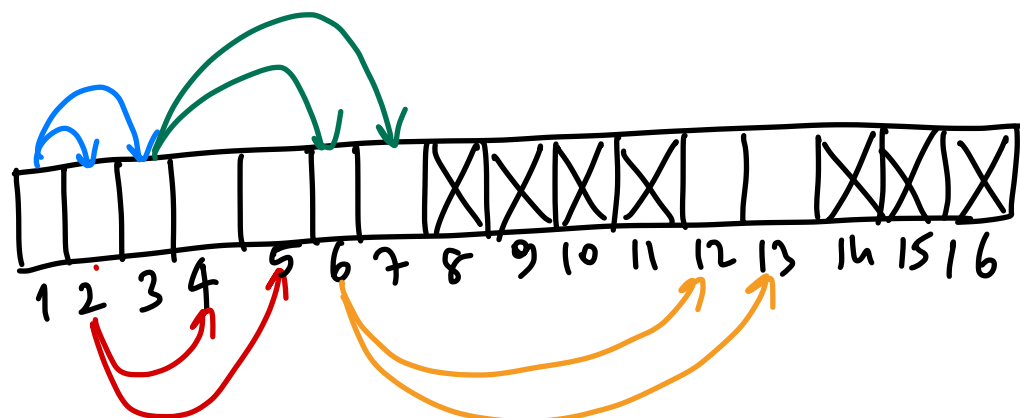
$$\lfloor 2.5 \rfloor = 2$$

$$5 = 2 * 2 + 1$$

$$\lfloor 3 \rfloor = 3$$

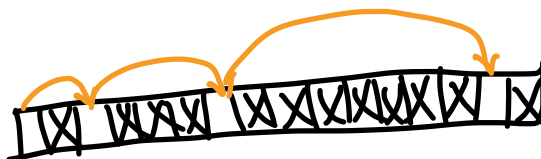
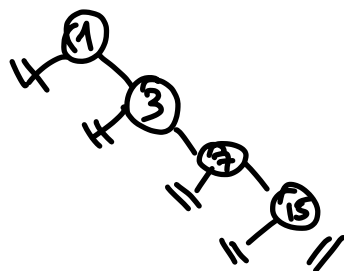
$$\text{parent} = \lfloor \text{child node} / 2 \rfloor$$

↑
floor function



$$6 = \lfloor 13 / 2 \rfloor = \lfloor 6.5 \rfloor$$

$$2 = \lfloor 4 / 2 \rfloor = \lfloor 2 \rfloor$$



Terminal nodes $\Rightarrow 4, 5, 7, 12, 13$

Path 1 $4 \rightarrow 2 \rightarrow 1$

Path 2 $5 \rightarrow 2 \rightarrow 1$

Path 3 $7 \rightarrow 3 \rightarrow 1$

Path 4 $12 \rightarrow 6 \rightarrow 3 \rightarrow 1$

Path 5 $13 \rightarrow 6 \rightarrow 3 \rightarrow 1$ ✓