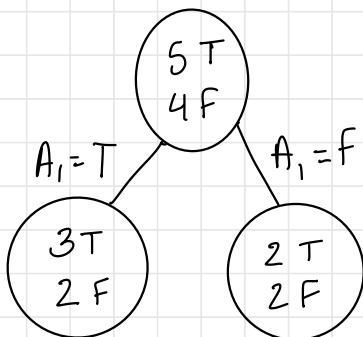


Assignment 3

(6a.) $P(+) = 5/9$; $P(-) = 4/9$

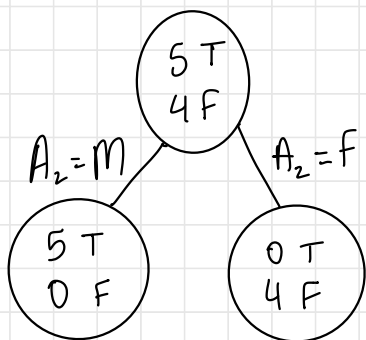
$$\text{Entropy} = -\frac{5}{9} \log_2 \left(\frac{5}{9} \right) + \frac{4}{9} \log_2 \left(\frac{4}{9} \right) = \boxed{-0.0489}$$

(6b.)



$$\begin{aligned} \text{Entropy: } & \frac{4}{9} \left[-\left(\frac{3}{5} \right) \log_2 \left(\frac{3}{5} \right) - \left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) \right] + \frac{5}{9} \left[-\left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) - \left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) \right] = \\ & 0.9871 \end{aligned}$$

$$\text{Gain}_{A_1} = -0.0489 - 0.9871 = -1.036$$



$$\begin{aligned} \text{Entropy: } & \frac{5}{9} \left[-\left(\frac{5}{5} \right) \log_2 \left(\frac{5}{5} \right) - \left(\frac{0}{5} \right) \log_2 \left(\frac{0}{5} \right) \right] + \frac{4}{9} \left[-\left(\frac{0}{4} \right) \log_2 \left(\frac{0}{4} \right) - \left(\frac{4}{4} \right) \log_2 \left(\frac{4}{4} \right) \right] = 0 \end{aligned}$$

$$\text{Gain}_{A_2} = -0.0489 - 0 = -0.0489$$

A3:

Sorted \rightarrow 2 3 4 5 6 7 8

Split Points $\rightarrow 1.5, 2.5, 3.5, 4.5, 5.5, 6.5, 7.5, 8.5$

Entropies \rightarrow

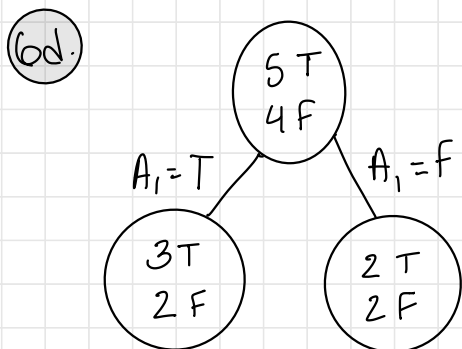
	<u>1.5</u>		<u>2.5</u>		<u>3.5</u>		<u>4.5</u>		<u>5.5</u>		<u>6.5</u>		<u>7.5</u>	
	\leq	$>$	\leq	$>$	\leq	$>$	\leq	$>$	\leq	$>$	\leq	$>$	\leq	$>$
+	0	5	1	4	2	3	3	2	3	2	3	2	4	1
-	0	4	0	4	0	4	0	4	2	2	3	1	3	1

	<u>8.5</u>	
	\leq	$>$
+	5	0
-	4	0

$-0.0489, .22, 0.17, 0.116, 0.107, 0.064, 0.039,$
 0

The best split for A3 is 8.5.

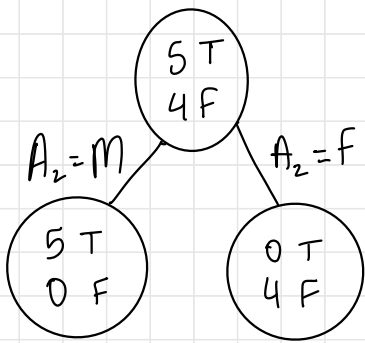
(6c.) A3 is the best split amongst the 3.



$$\text{Error} = 5/9 = 0.56$$

$$\begin{aligned} \text{Error} &= 5/9(2/5) + 4/9(2/4) \\ &= 0.44 \end{aligned}$$

$$\text{gain} = 0.11 \rightarrow A_1$$



$$\text{error} = \frac{4}{9} \left(\frac{5}{5} \right) + \frac{5}{9} \left(\frac{4}{4} \right)$$

$$= 1$$

$$\underline{\text{gain}} = -0.44 \rightarrow A_2$$

We should split based on A_1 .

(6e.) $\text{GINI} = 1 - (5/9)^2 - (4/9)^2 = 0.494$

$$\text{GINI} = \frac{5}{9} \left(1 - \left(\frac{3}{5} \right)^2 - \left(\frac{2}{5} \right)^2 \right) + \frac{4}{9} \left(1 - \left(\frac{2}{4} \right)^2 - \left(\frac{2}{4} \right)^2 \right)$$

$$= 0.489$$

$$\text{gain}_{A_1} = 0.005$$

$$\text{GINI} = \frac{5}{9} \left(1 - \left(\frac{5}{5} \right)^2 \right) + \frac{4}{9} \left(1 - 0^2 - \left(\frac{4}{4} \right)^2 \right)$$

$$= 0$$

$$\text{gain}_{A_2} = .494$$

so split at A_2

7a. The error on the training set is equal to the error on the testing set.

7b. Before split

$$\text{resubmission error} = 4/10$$

$$\text{generalization error} = 4.5/10$$

After split

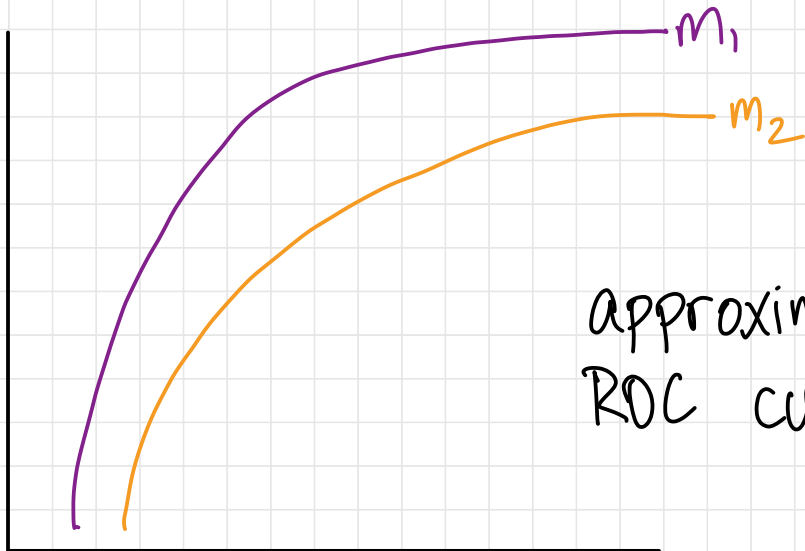
$$\text{resubmission error} =$$

$$\left(\frac{1}{1}\right)\left(\frac{1}{1}\right) + \left(\frac{1}{1}\right)\left(\frac{1}{1}\right) = 2$$

→ trim subtree

7c. Tree must be pruned/trimmed as in part (b).

8a.



approximate
ROC curves

8b. Precision = $\frac{TP}{TP+FP} = \frac{1}{1+4} = \frac{1}{5}$

Recall = $\frac{TP}{TP+FN} = \frac{1}{1+2} = \frac{1}{3}$

F-measure = $\frac{2(TP)}{2(TP) + FN + FP} = \frac{2}{2+2+4} = \frac{2}{8}$