

Mawaba Pascal Dao

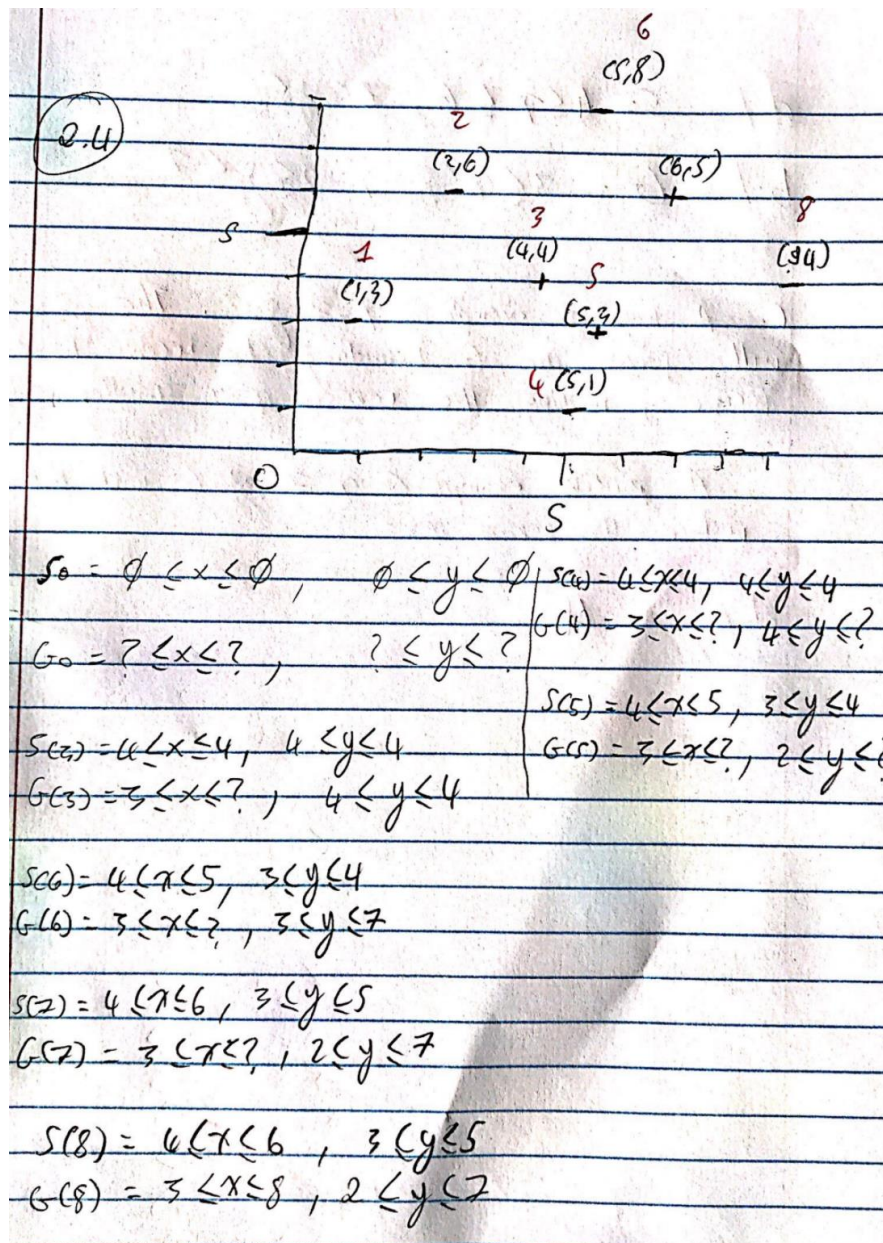
ML HW2

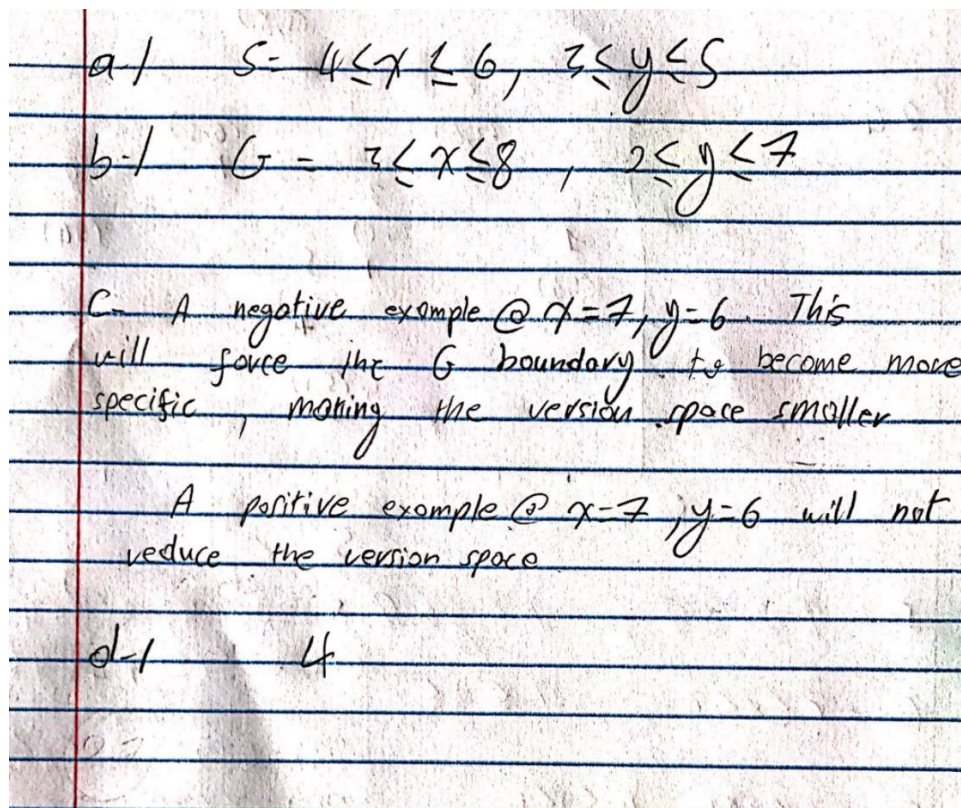
CSE5693

Dr. Chan

Feb 23, 2022

a.





b.

There cannot be a maximally specific consistent hypothesis for classifying real numbers because the interval between any 2 real numbers is infinite. In the example of $4.5 < x < 6.1$, if x was 5 a more specific consistent hypothesis could be $4.0 < x < 6.0$, and even more specific one could be $4.999 < x < 5.001$. The hypothesis can be made infinitely more specific using real numbers. We can modify the hypothesis representation to converge to a maximally specific consistent hypothesis by fixing the decimal point precision of the natural numbers, or by using integers instead of natural numbers.

c.

Ex. 11

$$Entropy(S) = -\left[\frac{3}{4} \log_2\left(\frac{3}{4}\right) + \left[-\frac{1}{4} \log_2\left(\frac{1}{4}\right)\right]\right]$$

$$= 0.75(-0.415) + (-0.25)(-2)$$

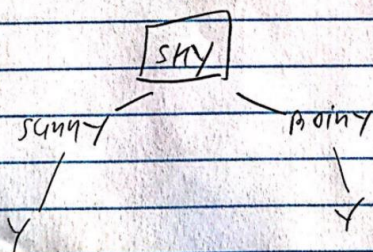
$$= 0.31125 + 0.5$$

$$Entropy(S) = 0.81125$$

Prut node selection

$$Gain(S, rky) = 0.81125 - \left[\frac{3}{4} (-1 \log_2(1) - \log_2(0)) + \frac{1}{4} (\log_2(0) - 1 \log_2(1)) \right]$$

We stop here because condition 2 is met



3.4 -b

The decision tree learned in this case was one with the sky attribute at the root, and with sky value ending up as the yes leaf node

and the rainy value ends up with the N leaf node. This is equivalent to the following general boundary hypothesis from the version space in Figure 2.3: <Sunny,?,?,?,?>. This hypothesis will classify on Sunny samples as positive

and all non-sunny samples as negative, the decision tree and version space from figure 2.3 therefore both output the same classification on this training set.

$$\begin{aligned}
 c/ \text{Entropy}(S) &= -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) \\
 &= -0.6(-0.7369) - 0.4(-1.3219) \\
 &= -0.6(-0.7369) - 0.4(-1.3219) \\
 &= 0.44214 + 0.52876 \\
 \text{Entropy}(S) &= 0.971 \\
 \text{Gain}(S, \text{Sky}) &= 0.971 - \left[\frac{4}{5} \left(-\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \right) + \frac{1}{5} (0) \right] \\
 &= 0.971 - \left[\frac{4}{5} (0.3127 + 0.5) \right] \\
 \text{Gain}(S, \text{Sky}) &= 0.3219 \\
 \text{Gain}(S, \text{AirTemp}) &= 0.971 - \left[\frac{4}{5} \left(-\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \right) + \frac{1}{5} (0) \right] \\
 \text{Gain}(S, \text{AirTemp}) &= 0.3219 \\
 \text{Gain}(S, \text{Humidity}) &= 0.971 - \left[\frac{2}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right) + \frac{3}{5} \left(-\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) \right) \right] \\
 &= 0.971 - \left(\frac{2}{5} (1) + \frac{3}{5} (0.3899 + 0.5287) \right) \\
 \text{Gain}(S, \text{Humidity}) &= 0.971 - 0.8802 \\
 \text{Gain}(S, \text{Humidity}) &= 0.0908
 \end{aligned}$$

$$G(S, \text{wind}) = 0.971 - \left[\frac{4}{5} \left(-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right) + \frac{1}{5} (0) \right]$$

$$G(S, \text{wind}) = 0.971$$

$$G(S, \text{water}) = 0.971 - \left[\frac{4}{5} \left(-\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \right) + \frac{1}{5} (0) \right]$$

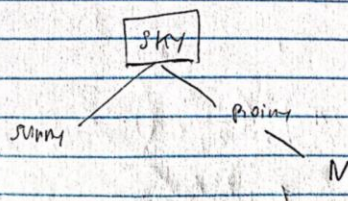
$$0.971 - \frac{4}{5} (1)$$

$$G(S, \text{water}) = 0.171$$

$$G(S, \text{Forecast}) = 0.971 - \left[\frac{3}{5} \left(-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right) + \frac{2}{5} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) \right]$$

$$= 0.971 - \left(\frac{3}{5} (0.9185) + \frac{2}{5} (1) \right)$$

$$G(S, \text{Forecast}) = 0.0199$$



$$G(\text{Sunny}, \text{Airtemp})$$

$$G(\text{Sunny}, \text{Humidity})$$

$$G(\text{Sunny}, \text{wind})$$

$$G(\text{Sunny}, \text{water})$$

$$G(\text{Sunny}, \text{Forecast})$$

$$\text{Entropy}(\text{Sunny}) = \frac{3}{4} \log_2 \left(\frac{3}{4} \right) + \frac{1}{4} \log_2 \left(\frac{1}{4} \right)$$

$$= 0.31127 + 0.5$$

$$\underline{E(\text{Sunny}) = 0.81127}$$

$$\text{Gain}(\text{Sunny}, \text{Pickup}) = 0.81127 - \left[1 \left(\frac{3}{4} \log_2 \left(\frac{3}{4} \right) + \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right) \right]$$

$$\text{Gain}(\text{Sunny}, \text{Pickup}) = 0.81127 - 0.81127 = 0$$

$$\begin{aligned} \text{Gain}(\text{Sunny}, H) &= 0.81127 - \left[\frac{2}{4} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) + \frac{2}{4} \log_2(0) \right] \\ &= 0.81127 - \frac{2}{4} \end{aligned}$$

$$\underline{\text{Gain}(\text{Sunny}, H) = 0.31127}$$

$$\underline{G(\text{Sunny}, \text{wind}) = 0.81127 - [0]}$$

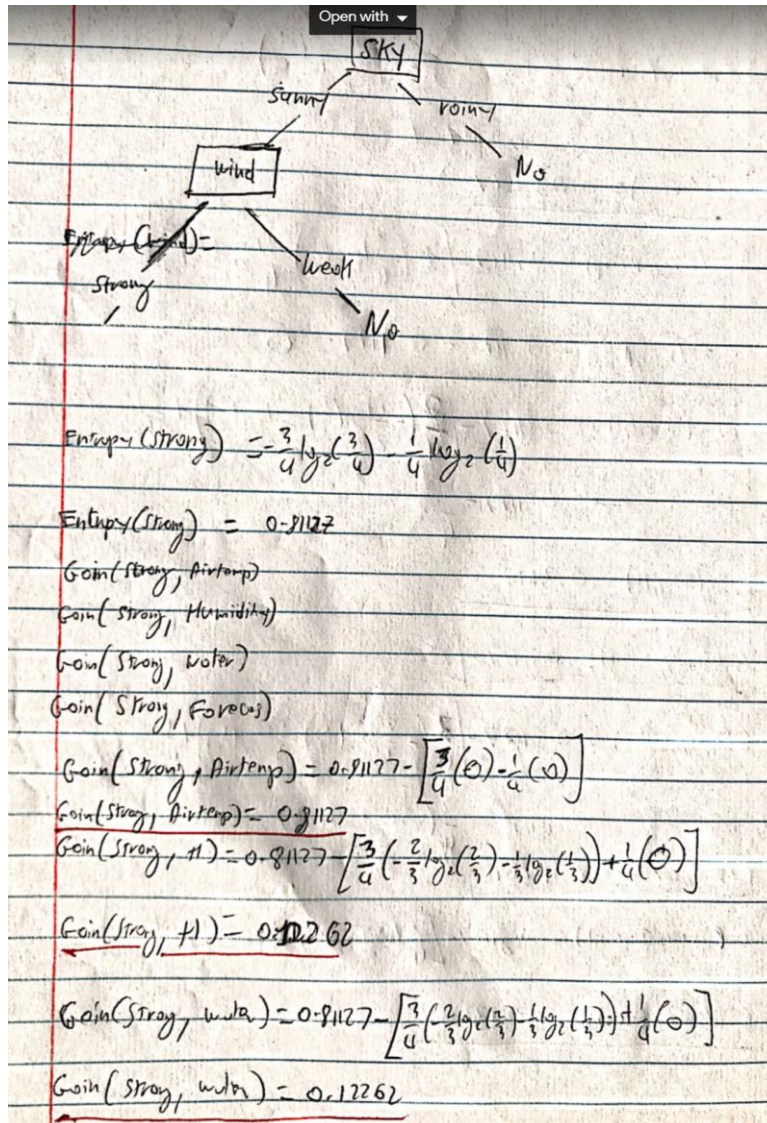
$$G(\text{Sunny}, \text{water}) = 0.81127 - \left[\frac{3}{4} \left(-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right) + \frac{1}{4} \log_2(0) \right]$$

$$0.81127 - \frac{3}{4} (0.3333 + 0.5283)$$

$$\underline{G(\text{Sunny}, \text{water}) = 0.12262}$$

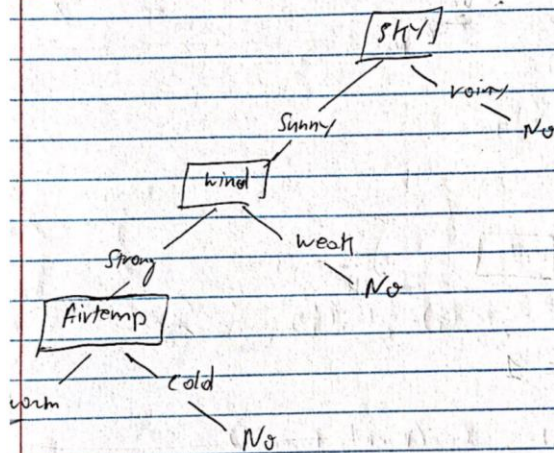
$$G(\text{Sunny}, \text{Forecast}) = 0.81127 - \left[\frac{3}{4} \left(-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right) + 0 \right]$$

$$\underline{\text{Gain}(\text{Sunny}, \text{Forecast}) = 0.12262}$$



$$\text{Gain}(\text{Strong}, \text{Forecast}) = 0.91127 - \left[\frac{2}{4} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) \right]$$

$$\text{Gain}(\text{Strong}, \text{Forecast}) = 0.31127$$



$$\text{Entropy}(\text{warm}) = \cancel{0.91127} \quad 0.81127$$

$$\text{Gain}(\text{warm}, H)$$

$$\text{Gain}(\text{warm}, \text{water})$$

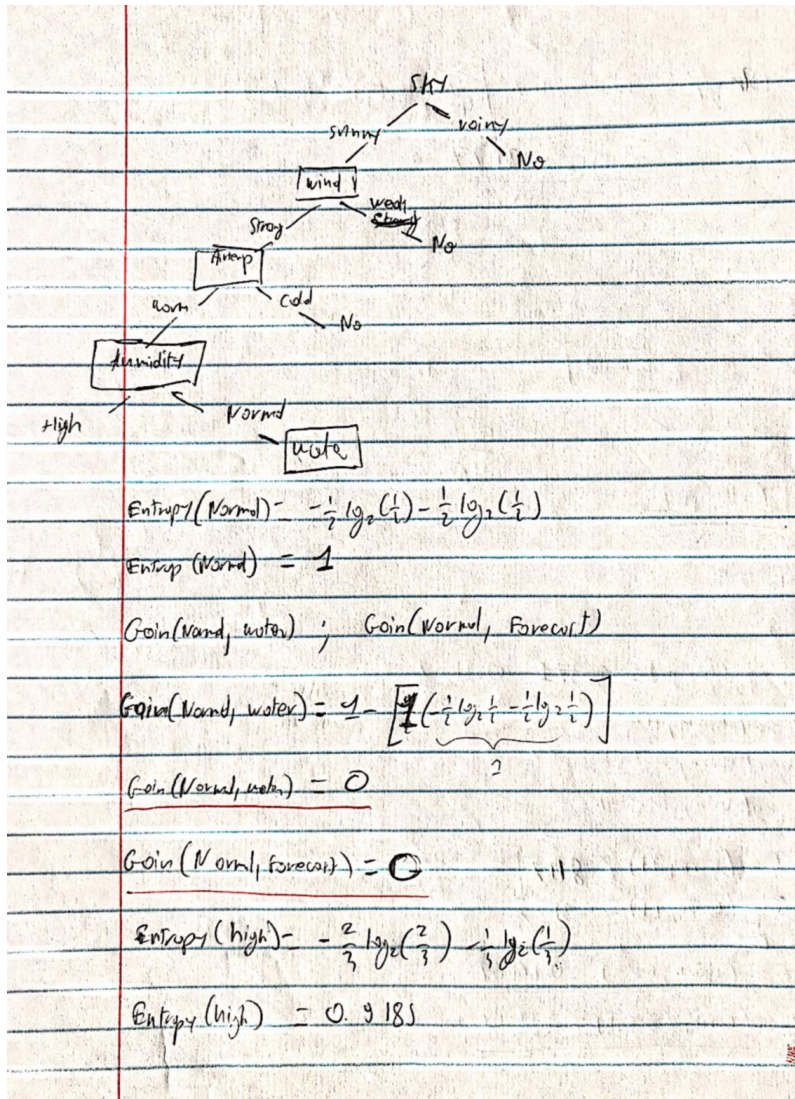
$$\text{Gain}(\text{warm}, \text{Forecast})$$

$$\text{Gain}(\text{warm}, H) = \cancel{0.91127}$$

$$\text{Gain}(\text{warm}, H) = 0.31127$$

$$\text{Gain}(\text{warm}, \text{water}) = 0.12262$$

$$\text{Gain}(\text{warm}, \text{Forecast}) = 0.12262$$



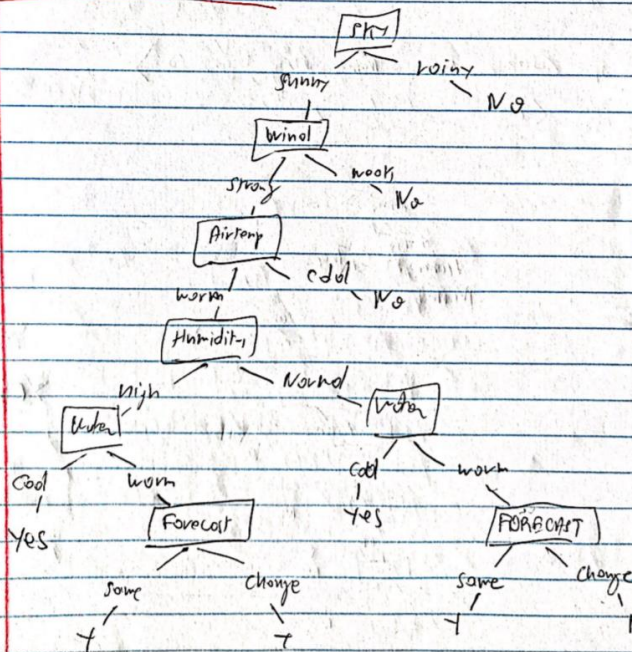
$$\text{Gain}(\text{high}, \text{wind}) = 0.9185 - \left[\frac{2}{3} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) + \frac{1}{3} (0) \right]$$

$$\text{Gain}(\text{high}, \text{wind}) = 0.25193$$

$$\text{Gain}(\text{high}, \text{Forecast}) = 0.9185 - \left[\frac{1}{3} (0) + \frac{2}{3} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) \right]$$

$$= 0.9185 - \frac{2}{3}$$

$$\text{Gain}(\text{high}, \text{Forecast}) = 0.25193$$



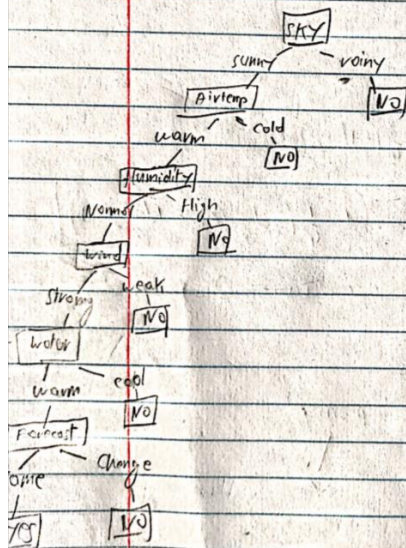
d-1

$S_1 = \langle \text{Sunny}, \text{warm}, \text{Normal}, \text{strong}, \text{warm}, \text{same} \rangle$

$G_1 = \langle ?, ?, ?, ?, ?, ? \rangle$

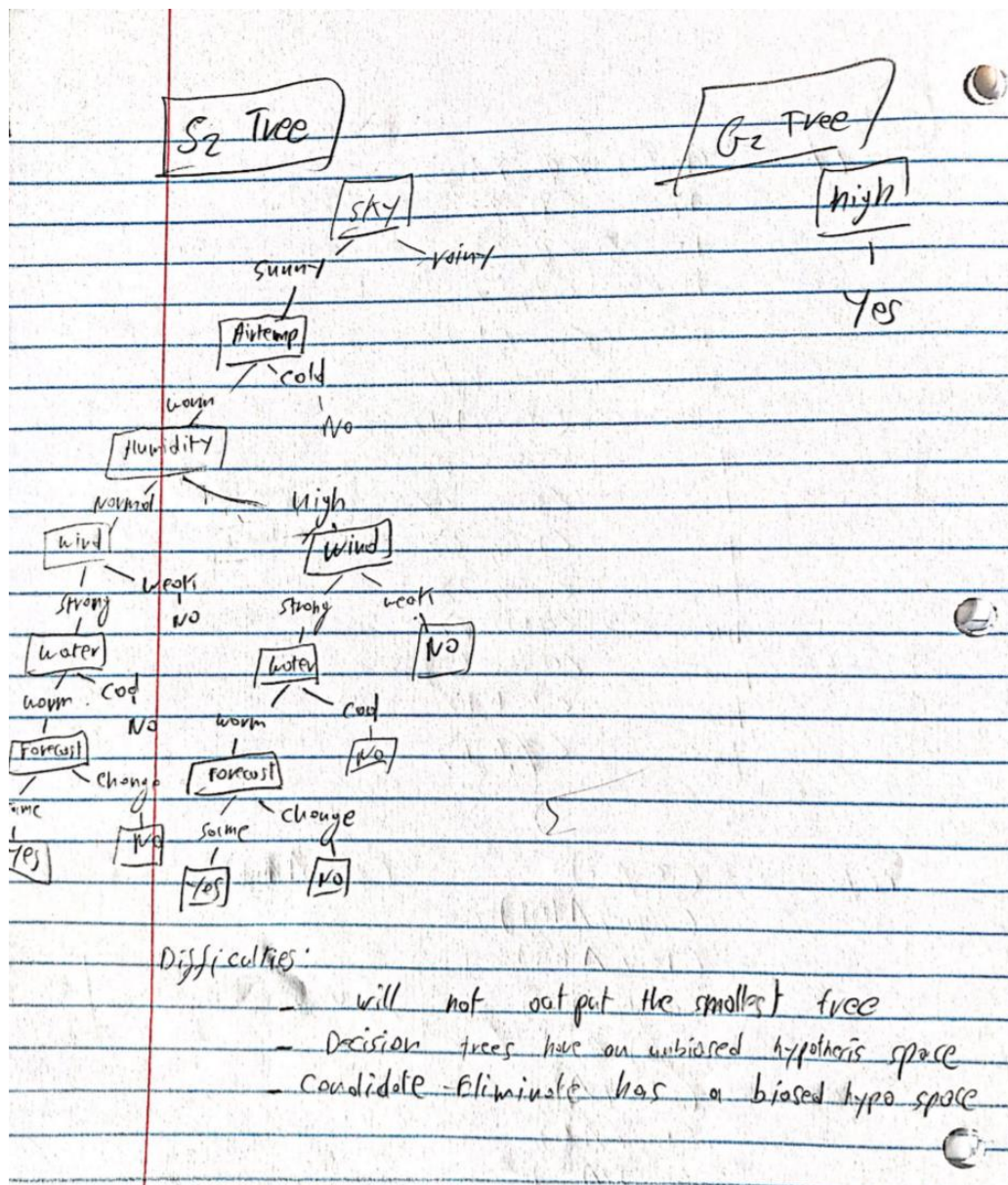
S_1 Tree:

G_1 Tree:



YES

root node



Candidate Eliminate does not use the information gain heuristic to select the next attribute and find the smallest tree. Candidate-Eliminate will output unnecessarily large trees in its specific boundary, each is a large # of syntactical equivalents.

d.

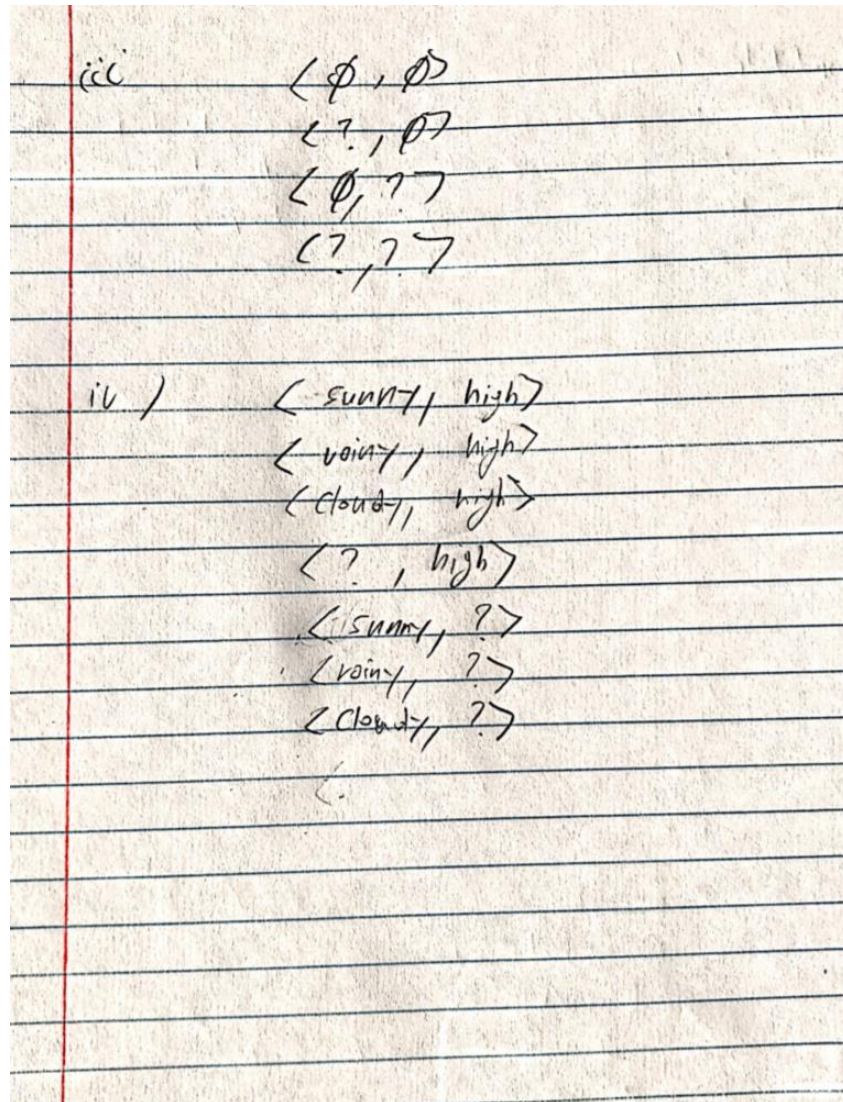
(d)

h_1	$\langle \emptyset, \emptyset \rangle$
h_2	$\langle ?, \emptyset \rangle$
h_3	$\langle \emptyset, ? \rangle$
h_4	$\langle ?, ? \rangle$
h_5	$\langle \text{sunny}, \text{high} \rangle$
h_6	$\langle \text{rainy}, \text{high} \rangle$
h_7	$\langle \text{cloudy}, \text{high} \rangle$
h_8	$\langle \text{sunny}, \emptyset \rangle$
h_9	$\langle \text{sunny}, ? \rangle$
h_{10}	$\langle \text{rainy}, \emptyset \rangle$
h_{11}	$\langle \text{rainy}, ? \rangle$
h_{12}	$\langle \text{cloudy}, \emptyset \rangle$
h_{13}	$\langle \text{cloudy}, ? \rangle$
h_{14}	$\langle ?, \text{high} \rangle$
h_{15}	$\langle \emptyset, \text{high} \rangle$

cc)

$\langle ? \wedge ? \rangle$	$\langle ? \wedge \text{high} \rangle$
$\langle \text{sunny} \wedge \text{high} \rangle$	
$\langle \text{rainy} \wedge \text{high} \rangle$	
$\langle \text{cloudy} \wedge \text{high} \rangle$	
$\langle \emptyset \rangle$	
$\langle \text{sunny} \wedge ? \rangle$	
$\langle \text{rainy} \wedge ? \rangle$	
$\langle \text{cloudy} \wedge ? \rangle$	

9 unique hypotheses



e.

Pruning the tree yielded better results on the test set as expected. The accuracy went from 76.4% to 88% after pruning.

Increasing the percentage of the corrupted data, testIrisNoisy, yields the same test accuracy from 2% to 20% during experimentation. This is likely due to a bug in testIrisNoisy that I did not have enough time to fix.

Accuracy plot

