

CSU 467 PS 5

1. a) $\frac{n-1}{n} = 1 - \frac{1}{n}$ since there is a $\frac{1}{n}$ chance that the 1^{st} observation is the j^{th} observation from the original sample, so $1 - \frac{1}{n}$ is the negation
- b) $\frac{n-1}{n} = 1 - \frac{1}{n}$ for the same reason as part a since the 2^{nd} and 1^{st} observations are independent. (with replacement)
- c) $(1 - \frac{1}{n})^n$ Since bootstrap draws observation with replacement, the chance of drawing the j^{th} observation from the original sample for any of the n draws in the bootstrap sample is $1 - \frac{1}{n}$. Since there's n observations in the bootstrap sample, and since each observation is independent since there is replacement, the chance of the bootstrap sample not containing the j^{th} observation is $(1 - \frac{1}{n})^n$.
- d) $1 - (1 - \frac{1}{5})^5 = 1 - (\frac{4}{5})^5 = 0.6723$
- e) $1 - (1 - \frac{1}{100})^{100} = 1 - (\frac{99}{100})^{100} = 0.6340$
- f) $1 - (1 - \frac{1}{10000})^{10000} = 1 - (\frac{9999}{10000})^{10000} = 0.6321$

When n is large, the percentage of overlap between two bootstrap samples is around $\boxed{1/3 = 33\%}$ of observations. Since each bootstrap sample will contain $\approx 2/3$ of the observations, with two samples there will be $\approx 1/3$ overlap.

$$\rightarrow \lim_{n \rightarrow \infty} (1 + \frac{\lambda}{n})^n = e^\lambda$$

$$\rightarrow \text{when } \lambda = -1, \text{ then } \lim_{n \rightarrow \infty} (1 - \frac{1}{n})^n = e^{-1} \approx \frac{1}{3}$$

which represents the chance that the j^{th} observation will not be in the bootstrap sample of n draws from the n total sample. Therefore, the bootstrap sample will contain around $2/3$ of the observations, which will amount to $\approx 1/3$ of overlap between 2 bootstrap samples that each contain $\approx 2/3$ of the observations.

\rightarrow bootstrap cannot be used like cross validation in training algorithms to estimate errors since it uses replacement to draw samples, so observations may overlap between validation and training sets which you can't have for estimating error.

$$x_{11}^2 + x_{21}^2 + 2x_{11}x_{21}$$

$$= (x_{11} + x_{21})(x_{11} + x_{21})$$

$$2. \quad n=2, p=2, \quad x_{11}=x_{12}, \quad x_{21}=x_{22}$$

$$y_1 + y_2 = 0 \quad x_{11} + x_{21} = 0 \quad x_{12} + x_{22} = 0 \quad \text{so estimate } \hat{\beta}_0 = 0$$

a) minimize:

$$\sum_{i=1}^2 (y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_{ij} x_{ij})^2 + \lambda \sum_{j=1}^p \hat{\beta}_{ij}^2$$

$$= [(y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12})^2 + (y_2 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{22})^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2)]$$

$$b) = y_1^2 - 2y_1 \hat{\beta}_1 x_{11} - 2y_1 \hat{\beta}_2 x_{12} + \hat{\beta}_1^2 x_{11}^2 + \hat{\beta}_2^2 x_{12}^2 + 2\hat{\beta}_1 \hat{\beta}_2 x_{11} x_{12}$$

$$+ y_2^2 - 2y_2 \hat{\beta}_1 x_{21} - 2y_2 \hat{\beta}_2 x_{22} + \hat{\beta}_1^2 x_{21}^2 + \hat{\beta}_2^2 x_{22}^2 + 2\hat{\beta}_1 \hat{\beta}_2 x_{21} x_{22}$$

$$+ \lambda \hat{\beta}_1^2 + \lambda \hat{\beta}_2^2$$

$$\frac{d}{d\hat{\beta}_1} = -2y_1 x_{11} + 2x_{11}^2 \hat{\beta}_1 + 2\hat{\beta}_2 x_{11} x_{12} - 2y_2 x_{21} + 2x_{21}^2 \hat{\beta}_1 + 2\hat{\beta}_2 x_{21} x_{22} + 2\lambda \hat{\beta}_1 = 0$$

$$\Rightarrow -y_1 x_{11} + x_{11}^2 \hat{\beta}_1 + \hat{\beta}_2 x_{11} x_{12} - y_2 x_{21} + x_{21}^2 \hat{\beta}_1 + \hat{\beta}_2 x_{21} x_{22} + \lambda \hat{\beta}_1 = 0$$

$$\text{with } x_{11} = x_{12}, \quad x_{21} = x_{22}$$

$$-y_1 x_{11} + x_{11}^2 \hat{\beta}_1 + \hat{\beta}_2 x_{11}^2 - y_2 x_{21} + x_{21}^2 \hat{\beta}_1 + \hat{\beta}_2 x_{21}^2 + 2\lambda \hat{\beta}_1$$

$$\lambda \hat{\beta}_1 + \hat{\beta}_1 (x_{11}^2 + x_{21}^2) + \hat{\beta}_2 (x_{11}^2 + x_{21}^2) = y_1 x_{11} + y_2 x_{21} \quad \text{and}$$

$$\Rightarrow \lambda \hat{\beta}_1 + \hat{\beta}_1 (x_{11}^2 + x_{21}^2) + \hat{\beta}_2 (x_{11}^2 + x_{21}^2) = y_1 x_{11} + y_2 x_{21} + 2\hat{\beta}_1 x_{11} x_{21} + 2\hat{\beta}_2 x_{11} x_{21}$$

$$\Rightarrow \lambda \hat{\beta}_1 = y_1 x_{11} + y_2 x_{21} + 2\hat{\beta}_1 x_{11} x_{21} + 2\hat{\beta}_2 x_{11} x_{21}$$

$$\frac{d}{d\hat{\beta}_2} = -2y_1 x_{12} + 2x_{12}^2 \hat{\beta}_2 + 2\hat{\beta}_1 x_{11} x_{12} - 2y_2 x_{22} + 2\hat{\beta}_2 x_{22}^2 + 2\hat{\beta}_1 x_{21} x_{22} + 2\lambda \hat{\beta}_2 = 0$$

$$\Rightarrow -y_1 x_{12} + x_{12}^2 \hat{\beta}_2 + \hat{\beta}_1 x_{11} x_{12} - y_2 x_{22} + \hat{\beta}_2 x_{22}^2 + \hat{\beta}_1 x_{21} x_{22} + \lambda \hat{\beta}_2 = 0$$

$$\text{with } x_{11} = x_{12}, \quad x_{21} = x_{22}$$

$$-y_1 x_{12} + x_{12}^2 \hat{\beta}_2 + \hat{\beta}_1 x_{11}^2, -y_2 x_{22} + \hat{\beta}_2 x_{22}^2 + \hat{\beta}_1 x_{21}^2 + \lambda \hat{\beta}_2 = 0$$

$$\Rightarrow \lambda \hat{\beta}_2 + \hat{\beta}_2 (x_{12}^2 + x_{22}^2) + \hat{\beta}_1 (x_{11}^2 + x_{21}^2) = y_1 x_{11} + y_2 x_{21}$$

$$\Rightarrow \lambda \hat{\beta}_2 = y_1 x_{11} + y_2 x_{21} + 2\hat{\beta}_1 x_{11} x_{21} + 2\hat{\beta}_2 x_{11} x_{21}$$

$$\text{so } \lambda \hat{\beta}_1 = \lambda \hat{\beta}_2$$

$$\text{so } \hat{\beta}_1 = \hat{\beta}_2$$

c) lasso : minimize $\text{RSS} + \lambda \sum_{j=1}^p |\beta_j|$

$$= (y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12})^2 + (y_2 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{22})^2 + \lambda(|\hat{\beta}_1| + |\hat{\beta}_2|)$$

d) lasso coefficients $\hat{\beta}_1$ and $\hat{\beta}_2$ are not unique in both $\hat{\beta}_2$ and $\hat{\beta}_1$ can either be ^{both} positive or ^{both} negative and produce the same solution.

$$\frac{d}{d\beta_1} \lambda |\hat{\beta}_1| = \lambda \frac{|\hat{\beta}_1|}{\hat{\beta}_1}$$

$$\frac{\lambda |\hat{\beta}_1|}{\hat{\beta}_2} = \frac{\lambda |\hat{\beta}_2|}{\hat{\beta}_2}$$

so solutions include:
or $\hat{\beta}_1$ and $\hat{\beta}_2$ are both positive
 $\hat{\beta}_1$ and $\hat{\beta}_2$ are both negative.

3. depth 1

PRESS COLOR pink : 5/8 icecream nonpink : 6/8 marshmallow
LIKES SINGING: yes : 5/8 icecream no : 6/8 marshmallow
LIKES DANCING: yes : 5/8 icecream no : 6/8 marshmallow
AGE: child : 5/8 icecream adult : 6/8 marshmallow

Select DRESS COLOR as 1st node

depth 2

with DRESS COLOR = pink

err rate

LIKES SINGING: yes : 4/4 icecream no : 3/4 marshmallow 0 + $\frac{1}{4}$
LIKES DANCING: yes : 3/4 icecream no : 2/4 marshmallow $\frac{1}{4} + \frac{1}{2}$
AGE: child : 3/4 icecream adult : 2/4 marshmallow $\frac{1}{4} + \frac{1}{2}$

Select LIKES SINGING as 2nd node

with DRESS COLOR = nonpink:

LIKES SINGING yes : 3/4 marshmallow no : 3/4 marshmallow $\frac{1}{4} + \frac{1}{4}$
LIKES DANCING yes : 2/4 marshmallow no : 4/4 marshmallow $\frac{1}{2} + 0$
AGE child : 2/4 icecream adult : 4/4 marshmallow $\frac{1}{2} + 0$

Select LIKES SINGING as 2nd node

depth 3

with PRESS COLOR = pink, LIKES SINGING = yes \Rightarrow classify to icecream

with DRESS COLOR = pink, LIKES SINGING = no :

LIKES DANCING: yes : 1/2 icecream no : 2/2 marshmallow $\frac{1}{2}$
AGE child : 1/2 icecream adult : 2/2 marshmallow $\frac{1}{2}$

Select LIKES DANCING as 3rd node

with DRESS COLOR = nonpink, LIKES SINGING = yes

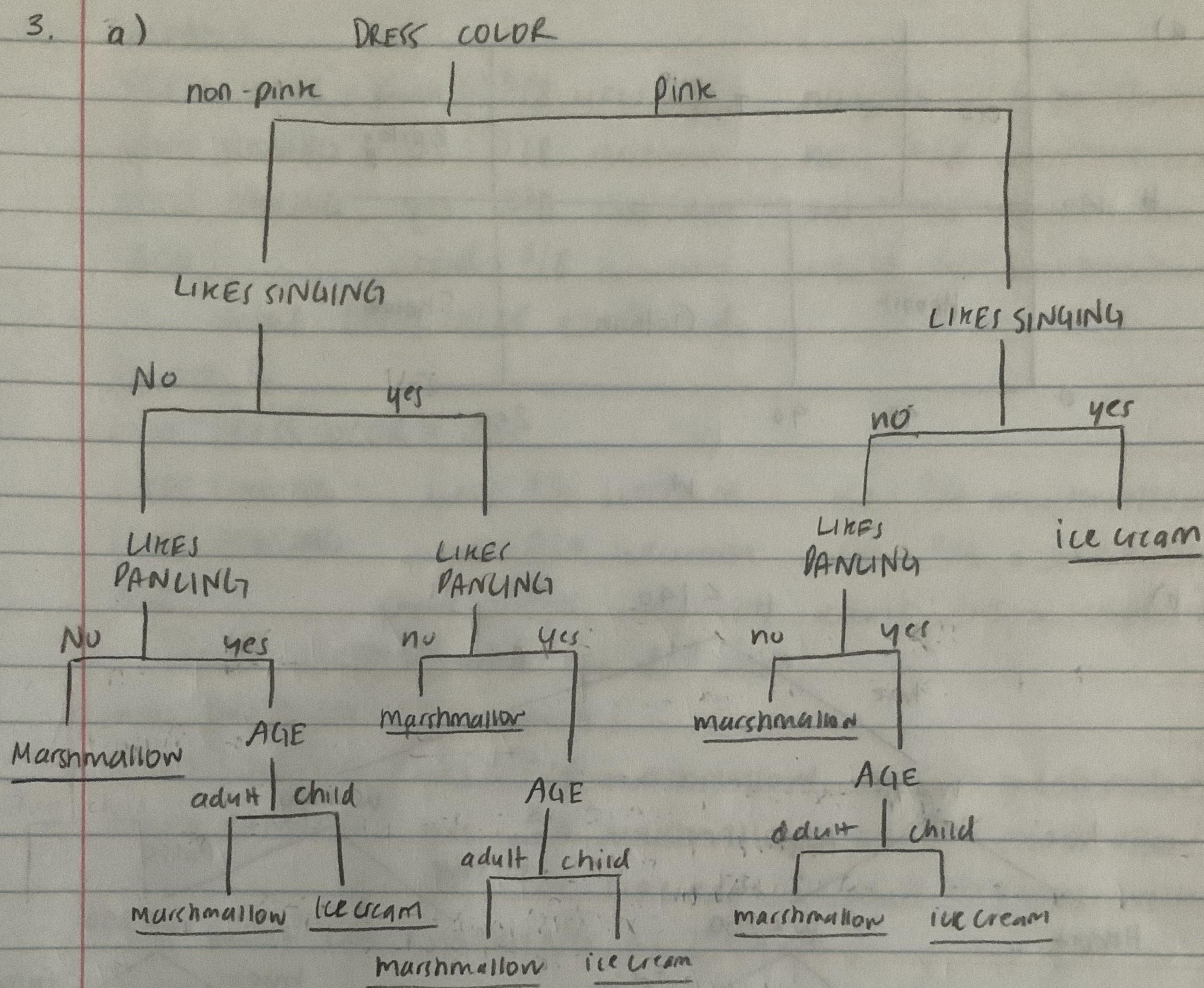
LIKES DANCING: yes : 1/2 icecream no : 2/2 marshmallow $\frac{1}{2}$
AGE child : 1/2 icecream adult : 2/2 marshmallow $\frac{1}{2}$

Select LIKES DANCING as 3rd node

with DRESS COLOR = nonpink, LIKES SINGING = no,

same as above, select LIKES DANCING as 3rd node

3. a)



b) IF DRESS COLOR = non pink & LIKES PANLING = yes & AGE = child

|| DRESS COLOR = pink & LIKES SINGING = yes

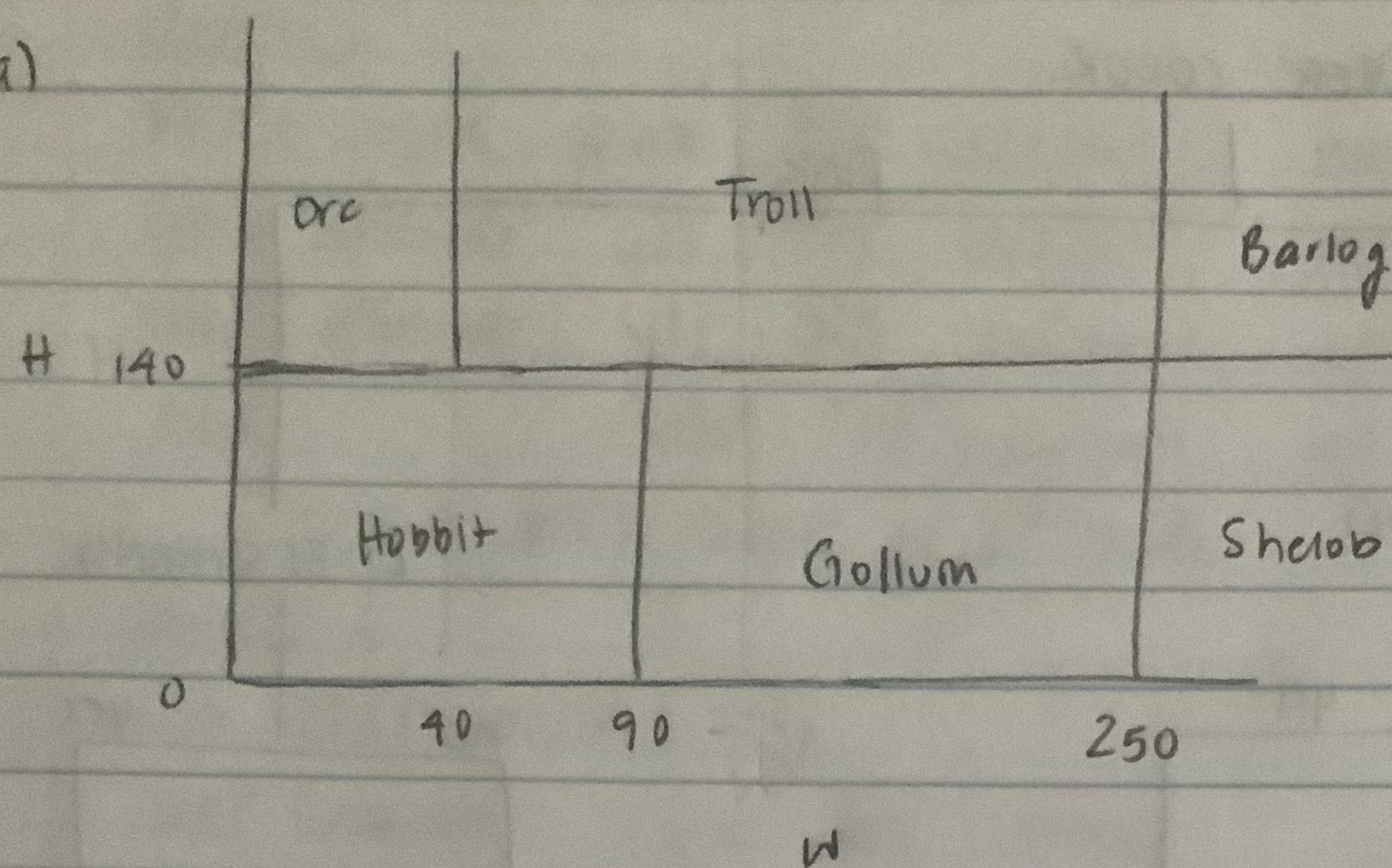
|| DRESS COLOR = pink & LIKES PANLING = yes & AGE = child)

THEN DESERT = ice cream

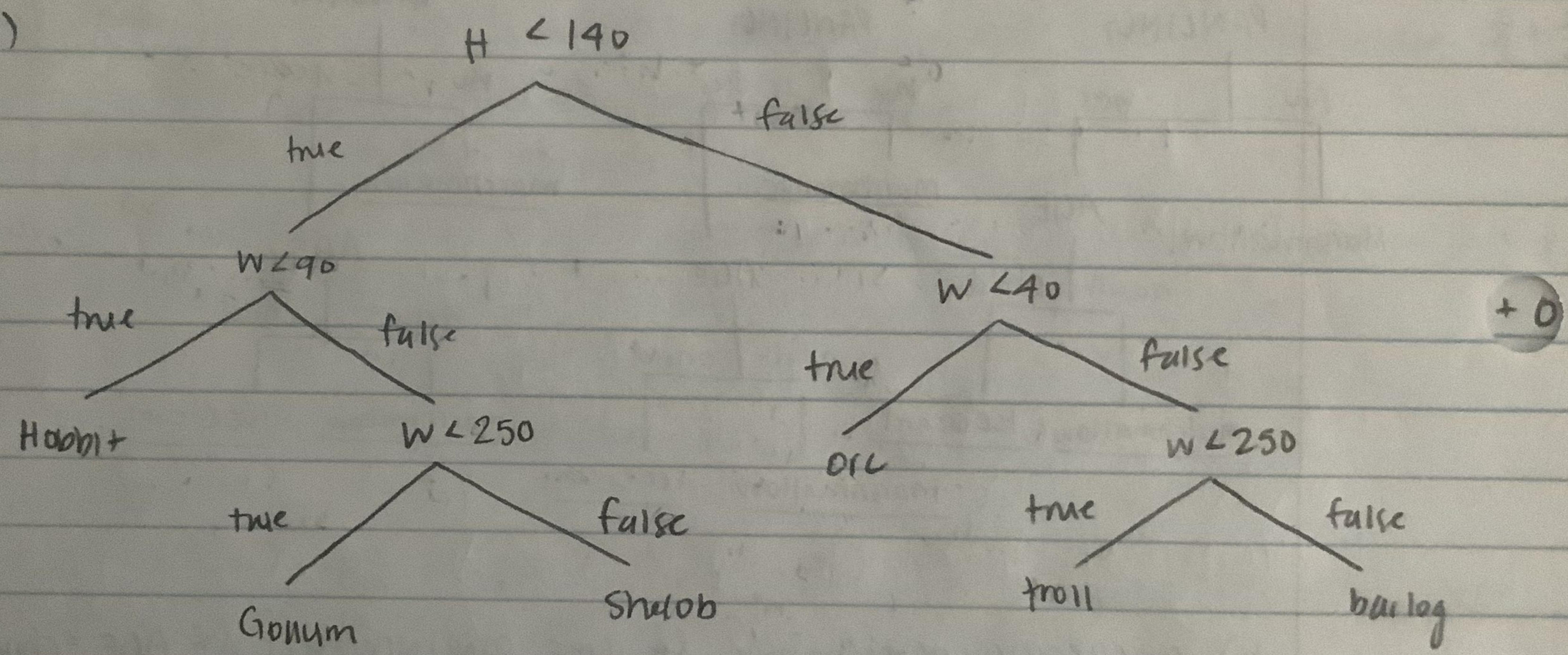
ELSE

THEN DESERT = marshmallow

4. a)



b)



5. a) use recursive ternary splitting

at each split, top down, make split such that

$$\{X \mid X_j < s_1\}, \{X \mid s_1 \leq X_j < s_2\} \text{ and } \{X \mid s_2 \leq X\}$$

leads to the greatest reduction in RSS where

$$\sum_{i: x_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2 + \sum_{i: x_i \in R_3(j,s)} (y_i - \hat{y}_{R_3})^2$$

$$\text{where } R_1(j,s) = \{X \mid X_j < s_1\}$$

$$R_2(j,s) = \{X \mid s_1 \leq X_j < s_2\}$$

$$R_3(j,s) = \{X \mid s_2 \leq X\}$$

then split one of the 3 resulting regions, and repeat.

b) until stopping criteria.

b) advantages:

- well suited for regions that clearly fit into 3 subregions
- more granularity at each node may lead to a tree of a lesser depth

disadvantages

- more computationally difficult to find optimal 3 regions at thresholds s_1 and s_2
- assumes the form that regions split into 3 subregions which won't work well for regions that don't fit assumption

c) You observe the form of the data through graphs

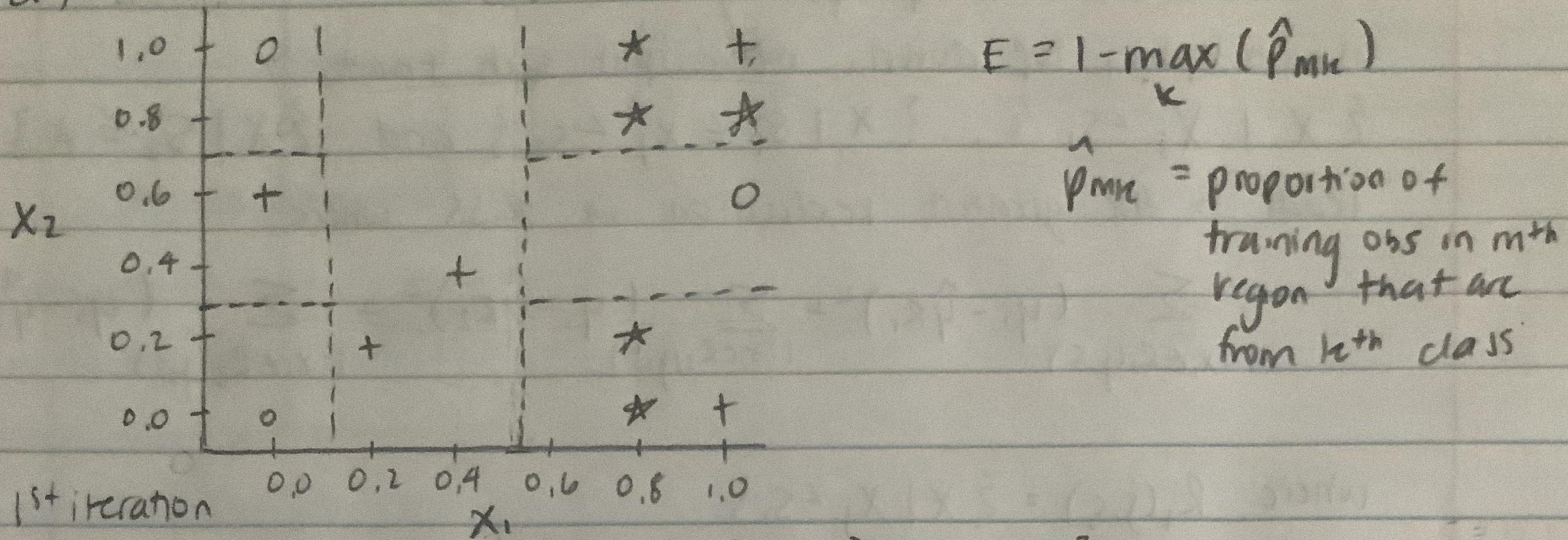
to see if regions seem as if they would be better suited to two way splits or three way splits.

Or you can construct both trees and see which produces a model with lower error rates.

Since binary splitting is more computationally advantageous.

Pick binary splitting if you need better computation efficiency.

d)



$$E = 1 - \max(\hat{p}_{mk})$$

\hat{p}_{mk} = proportion of
training obs in m^{th}
region that are
from k^{th} class

$$X_1: S_1 = 0.1, S_2 = 0.5 \quad E = (1 - \frac{2}{3}) + 0 + (1 - \frac{5}{8}) = \frac{1}{3} + \frac{3}{8} = 0.70833$$

$$X_1: S_1 = 0.5, S_2 = 0.9 \quad E = (1 - \frac{3}{5}) + 0 + (1 - \frac{1}{2}) = \frac{2}{5} + 0.5 = 0.9$$

so split X_1 at $S_1 = 0.1, S_2 = 0.5$

2nd iteration: splitting leftmost region, $X_1 \leq 0.1$

$$X_2: S_1 = 0.7, S_2 = 0.3, E = 0$$

3rd iteration: splitting rightmost region $X_1 > 0.5$

$$X_2: S_1 = 0.7, S_2 = 0.3 \quad E = (1 - \frac{3}{4}) + 0 + (1 - \frac{3}{4}) = \frac{1}{4} + \frac{1}{4} = 0.5$$

so result:

if $X_1 \leq 0.1$:

 if $X_2 \leq 0.3$ classify to "0"

 if $0.3 \leq X_2 \leq 0.7$ classify to "+"

 if $X_2 \geq 0.7$ classify to "0"

 if $0.1 \leq X_1 \leq 0.5$ classify to "+"

 if $0.5 \leq X_1$:

 if $X_2 \leq 0.3$ classify to "x"

 if $0.3 \leq X_2 \leq 0.7$ classify to "o"

 if $0.7 \leq X_2$ classify to "x"

$$6. \Pr(Y=1|X) : 0.09, 0.13, 0.19, 0.19, 0.53, 0.59, 0.6, 0.63, 0.71$$

$\downarrow \quad \downarrow \quad \downarrow$

majority polling: $Y=0 \quad Y=0 \quad Y=0 \quad Y=0 \quad Y=1 \quad Y=1 \quad Y=1 \quad Y=1 \quad Y=1$
since # $Y=1$'s = 6, final classification would be $\boxed{Y=1}$
which is the majority since $n=9$

average probability: $\frac{0.09+0.13+0.19+0.19+0.53+0.59+0.6+0.63+0.71}{9}$
= 0.40667

since average < 0.5, would classify as not $Y=1$,
so $\boxed{Y=0}$