

CSCC11 Assignment 2

Jefferson Li (jeffersonli.li@mail.utoronto.ca)

October 18, 2020

Question 1.

our goal is to maximize $P(f = 3)$

$$P(f = 3) = P(S + D = 3) \quad (1)$$

$$= P(S = s, D = 3 - s) = P(S = s)P(D = 3 - s) \quad \text{as they are independent} \quad (2)$$

$$P(S = s)P(D = 3 - s) = \left(\frac{1}{\sqrt{2\pi}3} \exp\left(-\frac{1}{2}\left(\frac{s-0}{9}\right)^2\right)\right) \left(\frac{1}{\sqrt{2\pi}2} \exp\left(-\frac{1}{2}\left(\frac{(3-s)-0}{4}\right)^2\right)\right) \quad \text{Replace with PDFs} \quad (3)$$

$$\ln(P(f = 3)) = \ln\left(\frac{1}{\sqrt{2\pi}3}\right) + \frac{-1}{2}\left(\frac{s-0}{9}\right)^2 + \ln\left(\frac{1}{\sqrt{2\pi}2}\right) + \frac{-1}{2}\left(\frac{(3-s)-0}{4}\right)^2 \quad \text{take ln for easier differentiation} \quad (4)$$

$$\frac{\partial \ln(P(f = 3))}{\partial s} = -1\left(\frac{s}{9}\right) - 1\left(\frac{3-s}{4}\right) * -1 \quad \text{differentiate the probability function} \quad (5)$$

$$\frac{\partial \ln(P(f = 3))}{\partial s} = \frac{-s}{9} + \frac{3}{4} - \frac{s}{4} = \frac{-13s}{36} + \frac{3}{4} \quad (6)$$

$$\frac{\partial \ln(P(f = 3))}{\partial s} = 0 \implies s = \frac{27}{13} \quad \text{find the most probable } s \quad (7)$$

$$\frac{\partial^2 \ln(P(f = 3))}{\partial^2 s} = \frac{-13}{36} < 0 \quad \text{2nd derivative check} \quad (8)$$

Since 2nd derivative is negative, $s = \frac{27}{13}$ is a global maximum and therefore the most probable starting location

Question 2.

- a) Using Uni-variate splits, we need to consider 4 splits. As there are 4 features, and each feature has 2 unique values. So each feature has 1 split, and therefore 4 splits in total

- b)

$$H(D_0) = P_{true} \log_2\left(\frac{1}{P_{true}}\right) + P_{false} \log_2\left(\frac{1}{P_{false}}\right) = \frac{7}{16} \log_2\left(\frac{16}{7}\right) + \frac{9}{16} \log_2\left(\frac{16}{9}\right) \approx 0.98869$$

- For the color split
let t_0 be -1 if color is Red, and 1 if color is Green

$$IG(D_0, t_0) = H(D_0) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (9)$$

$$= H(D_0) - \frac{8}{16} \left(\frac{5}{8} \log_2 \left(\frac{8}{6} \right) + \frac{3}{8} \log_2 \left(\frac{8}{3} \right) \right) - \frac{8}{16} \left(\frac{2}{8} \log_2 \left(\frac{8}{2} \right) + \frac{6}{8} \log_2 \left(\frac{8}{6} \right) \right) \quad (10)$$

$$\approx 0.105833 \quad (11)$$

- For the size split
let t_0 be -1 if size is Small, and 1 if size is Large

$$IG(D_0, t_0) = H(D_0) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (12)$$

$$= H(D_0) - \frac{8}{16} \left(\frac{5}{8} \log_2 \left(\frac{8}{6} \right) + \frac{3}{8} \log_2 \left(\frac{8}{3} \right) \right) - \frac{8}{16} \left(\frac{2}{8} \log_2 \left(\frac{8}{2} \right) + \frac{6}{8} \log_2 \left(\frac{8}{6} \right) \right) \quad (13)$$

$$\approx 0.105833 \quad (14)$$

- For the shape split
let t_0 be -1 if shape is Sphere, and 1 if shape is Cube

$$IG(D_0, t_0) = H(D_0) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (15)$$

$$= H(D_0) - \frac{8}{16} \left(\frac{5}{8} \log_2 \left(\frac{8}{6} \right) + \frac{3}{8} \log_2 \left(\frac{8}{3} \right) \right) - \frac{8}{16} \left(\frac{2}{8} \log_2 \left(\frac{8}{2} \right) + \frac{6}{8} \log_2 \left(\frac{8}{6} \right) \right) \quad (16)$$

$$\approx 0.105833 \quad (17)$$

- For the texture split
let t_0 be -1 if texture is Soft, and 1 if texture is Hard

$$IG(D_0, t_0) = H(D_0) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (18)$$

$$= H(D_0) - \frac{8}{16} \left(\frac{5}{8} \log_2 \left(\frac{8}{6} \right) + \frac{3}{8} \log_2 \left(\frac{8}{3} \right) \right) - \frac{8}{16} \left(\frac{2}{8} \log_2 \left(\frac{8}{2} \right) + \frac{6}{8} \log_2 \left(\frac{8}{6} \right) \right) \quad (19)$$

$$\approx 0.105833 \quad (20)$$

- c) NOTE : please use the part d) diagram to follow the node and level

Since all splits result in the same information gain, we can choose any one. Lets choose to split by color

- level 1 left node

$$H(D_1) = P_{true} \log_2 \left(\frac{1}{P_{true}} \right) + P_{False} \log_2 \left(\frac{1}{P_{False}} \right) = \frac{5}{8} \log_2 \left(\frac{8}{5} \right) + \frac{3}{8} \log_2 \left(\frac{8}{3} \right) \approx 0.954434$$

- * For the size split
let t_1 be -1 if size is Small, and 1 if size is Large

$$IG(D_1, t_1) = H(D_1) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (21)$$

$$= H(D_1) - \frac{1}{2} \left(\frac{1}{1} \log_2 \left(\frac{1}{1} \right) + \frac{0}{0} \log_2(0) \right) - \frac{1}{2} \left(\frac{1}{4} \log_2 \left(\frac{4}{1} \right) + \frac{3}{4} \log_2 \left(\frac{4}{3} \right) \right) \quad (22)$$

$$\approx 0.68 \quad (23)$$

* For the shape split

let t_1 be -1 if shape is Sphere, and 1 if shape is Cube

$$IG(D_1, t_1) = H(D_1) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (24)$$

$$= H(D_1) - \frac{1}{2} \left(\frac{3}{4} \log_2 \left(\frac{4}{3} \right) + \frac{1}{4} \log_2 \left(\frac{4}{1} \right) \right) - \frac{1}{2} \left(\frac{2}{4} \log_2 \left(\frac{4}{2} \right) + \frac{2}{4} \log_2 \left(\frac{4}{2} \right) \right) \quad (25)$$

$$\approx 0.36 \quad (26)$$

* For the texture split

let t_1 be -1 if texture is Soft, and 1 if texture is Hard

$$IG(D_1, t_1) = H(D_1) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (27)$$

$$= H(D_1) - \frac{1}{2} \left(\frac{3}{4} \log_2 \left(\frac{4}{3} \right) + \frac{1}{4} \log_2 \left(\frac{4}{1} \right) \right) - \frac{1}{2} \left(\frac{2}{4} \log_2 \left(\frac{4}{2} \right) + \frac{2}{4} \log_2 \left(\frac{4}{2} \right) \right) \quad (28)$$

$$\approx 0.36 \quad (29)$$

Using the largest IG, the split function we will use for level 1 left node is splitting by size
– level 1 right node

$$H(D_2) = P_{true} \log_2 \left(\frac{1}{P_{true}} \right) + P_{false} \log_2 \left(\frac{1}{P_{false}} \right) = \frac{2}{8} \log_2 \left(\frac{8}{2} \right) + \frac{6}{8} \log_2 \left(\frac{8}{6} \right) \approx .81127$$

* For the size split

let t_2 be -1 if size is Small, and 1 if size is Large

$$IG(D_2, t_2) = H(D_2) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (30)$$

$$= H(D_2) - \frac{1}{2} \left(\frac{1}{4} \log_2 \left(\frac{4}{1} \right) + \frac{3}{4} \log_2 \left(\frac{4}{3} \right) \right) - \frac{1}{2} \left(\frac{1}{4} \log_2 \left(\frac{4}{1} \right) + \frac{3}{4} \log_2 \left(\frac{4}{3} \right) \right) \quad (31)$$

$$\approx 0.311 \quad (32)$$

* For the shape split

let t_2 be -1 if shape is Sphere, and 1 if shape is Cube

$$IG(D_2, t_2) = H(D_2) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (33)$$

$$= H(D_2) - \frac{1}{2} \left(\frac{2}{4} \log_2 \left(\frac{4}{2} \right) + \frac{2}{4} \log_2 \left(\frac{4}{2} \right) \right) - \frac{1}{2} (0) \quad (34)$$

$$\approx 0.56127 \quad (35)$$

* For the texture split

let t_2 be -1 if texture is Soft, and 1 if texture is Hard

$$IG(D_2, t_2) = H(D_2) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (36)$$

$$= H(D_2) - \frac{1}{2} \left(\frac{2}{4} \log_2 \left(\frac{4}{2} \right) + \frac{2}{4} \log_2 \left(\frac{4}{2} \right) \right) - \frac{1}{2} (0) \quad (37)$$

$$\approx 0.56127 \quad (38)$$

Using the largest IG, the split function we will use for level 1 right node is splitting by shape.

– level 2 left most node

* Since every data in this partition is True, we can terminate splitting and have a leaf node for *True*.

– level 2 left middle node

$$H(D_3) = P_{true} \log_2\left(\frac{1}{P_{true}}\right) + P_{false} \log_2\left(\frac{1}{P_{false}}\right) = \frac{1}{4} \log_2(4) + \frac{3}{4} \log_2\left(\frac{4}{3}\right) \approx 0.81127$$

* For the shape split

let t_3 be -1 if shape is Sphere, and 1 if shape is Cube

$$IG(D_3, t_3) = H(D_3) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (39)$$

$$= H(D_3) - \frac{1}{2} \left(\frac{1}{2} \log_2\left(\frac{2}{1}\right) + \frac{1}{2} \log_2\left(\frac{2}{1}\right) \right) - \frac{1}{2}(0) \quad (40)$$

$$\approx 0.56127 \quad (41)$$

* For the texture split

let t_3 be -1 if texture is Soft, and 1 if texture is Hard

$$IG(D_3, t_3) = H(D_3) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (42)$$

$$= H(D_3) - \frac{1}{2} \left(\frac{1}{2} \log_2\left(\frac{2}{1}\right) + \frac{1}{2} \log_2\left(\frac{2}{1}\right) \right) - \frac{1}{2}(0) \quad (43)$$

$$\approx 0.56127 \quad (44)$$

Using largest IG, we will split on shape

– level 2 right middle node

$$H(D_4) = P_{true} \log_2\left(\frac{1}{P_{true}}\right) + P_{false} \log_2\left(\frac{1}{P_{false}}\right) = \frac{1}{2} \log_2(2) + \frac{1}{2} \log_2\left(\frac{2}{1}\right) = 1$$

* For the size split

let t_4 be -1 if size is Small, and 1 if size is Large

$$IG(D_4, t_4) = H(D_4) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (45)$$

$$= H(D_4) - \frac{1}{2} \left(\frac{1}{2} \log_2\left(\frac{2}{1}\right) + \frac{1}{2} \log_2\left(\frac{2}{1}\right) \right) - \frac{1}{2} \left(\frac{1}{2} \log_2\left(\frac{2}{1}\right) + \frac{1}{2} \log_2\left(\frac{2}{1}\right) \right) \quad (46)$$

$$= 1 \quad (47)$$

* For the texture split

let t_4 be -1 if texture is Soft, and 1 if texture is Hard

$$IG(D_4, t_4) = H(D_4) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (48)$$

$$= H(D_4) - \frac{1}{2} \left(\frac{1}{1} \log_2\left(\frac{1}{1}\right) + \frac{0}{0} \log_2\left(\frac{0}{0}\right) \right) - \frac{1}{2} \left(\frac{0}{0} \log_2\left(\frac{0}{0}\right) + \frac{1}{1} \log_2\left(\frac{1}{1}\right) \right) \quad (49)$$

$$\approx 0.68627 \quad (50)$$

Using largest IG, we will split on texture

– level 2 right most node

* Since every data in this partition is False, we can terminate splitting and have a leaf node for *False*.

– level 3 left most node

$$H(D_5) = P_{true} \log_2\left(\frac{1}{P_{true}}\right) + P_{False} \log_2\left(\frac{1}{P_{False}}\right) = \frac{1}{2} \log_2(2) + \frac{1}{2} \log_2\left(\frac{2}{1}\right) = 1$$

* For the texture split

let t_5 be -1 if texture is Soft, and 1 if texture is Hard

$$IG(D_5, t_5) = H(D_5) - \frac{N_L}{N_j} H(D_L) - \frac{N_R}{N_j} H(D_R) \quad (51)$$

$$= H(D_5) - \frac{1}{2} \left(\frac{1}{1} \log_2\left(\frac{1}{1}\right) + \frac{0}{0} \log_2\left(\frac{0}{0}\right) \right) - \frac{1}{2} \left(\frac{0}{0} \log_2\left(\frac{0}{0}\right) + \frac{1}{1} \log_2\left(\frac{1}{1}\right) \right) \quad (52)$$

$$= 1 \quad (53)$$

Using largest IG, we will split on texture

– level 3 left middle node

* Since every data in this partition is False, we can terminate splitting and have a leaf node for *False*.

– level 3 right middle node

* Since every data in this partition is True, we can terminate splitting and have a leaf node for *True*.

– level 3 right most node

* Since every data in this partition is False, we can terminate splitting and have a leaf node for *False*.

– level 4 left node

* Since every data in this partition is True, we can terminate splitting and have a leaf node for *True*.

– level 4 right node

* Since every data in this partition is False, we can terminate splitting and have a leaf node for *False*.

- d)

