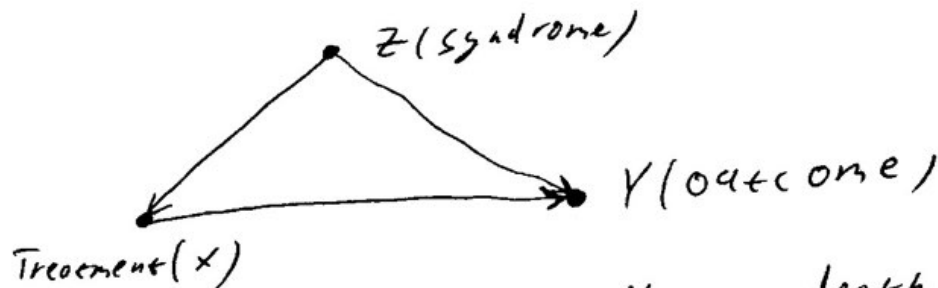


# Exercis 6

## 40211-1105.N

E1



$Z = z_1$  - presence of syndrome  
 $Z = z_0$  - absence of syndrome  
 $Y = y_1$  - death  
 $Y = y_0$  - survival  
 $X = x_1$  - take drug  
 $X = x_0$  - don't take drug

$$\begin{aligned}
 P(Z = z_1) &= r \\
 P(Z = z_0) &= 1 - r \\
 P(Y = y_1 | Z = z_0, X = x_1) &= P_2 \\
 P(Y = y_1 | Z = z_0, X = x_0) &= P_1 \\
 P(Y = y_1 | Z = z_1, X = x_0) &= P_3 \\
 P(Y = y_1 | Z = z_1, X = x_1) &= P_4
 \end{aligned}
 \left. \vphantom{\begin{aligned} P(Y = y_1 | Z = z_0, X = x_1) \\ P(Y = y_1 | Z = z_0, X = x_0) \\ P(Y = y_1 | Z = z_1, X = x_0) \\ P(Y = y_1 | Z = z_1, X = x_1) \end{aligned}} \right\} \begin{array}{l} \text{for case } Y = y_0 \\ \text{is to minus } 1 \\ \underline{\underline{1 - P(\dots)}} \end{array}$$

$$\begin{aligned}
 q_1 &= P(X = x_1 | Z = z_0) \\
 q_2 &= P(X = x_1 | Z = z_1)
 \end{aligned}
 \left. \vphantom{\begin{aligned} q_1 \\ q_2 \end{aligned}} \right\} \begin{array}{l} \text{for case } X = x_0 \\ \text{is to minus } 1 \\ \underline{\underline{1 - P(\dots)}} \end{array}$$

$$P(X, Y, Z) = P(Z) \cdot \underset{\substack{\uparrow \\ \text{parent}}}{P(X|Z)} \cdot \underset{\substack{\uparrow \\ \text{parent}}}{P(Y|X, Z)}$$

Let's find for all different values of  $X, Y, Z$ .

$$\begin{aligned}
 P(X = x_0, Y = y_0, Z = z_0) &= P(Z = z_0) \cdot P(X = x_0 | Z = z_0) \cdot P(Y = y_0 | X = x_0, Z = z_0) \\
 &= (1 - r)(1 - q_1)(1 - P_1)
 \end{aligned}$$

$$P(X=x_0, Y=y_1, Z=z_0) = (1-r)(1-q_1)P_1$$

$$P(X=x_1, Y=y_0, Z=z_0) = P(Z=z_0)P(X=x_1|Z=z_0)P(Y=y_0|X=x_1, Z=z_0)$$

$$= (1-r)q_1(1-P_2)$$

$$P(X=x_1, Y=y_1, Z=z_0) = (1-r)q_1P_2$$

$$P(X=x_0, Y=y_0, Z=z_1) = r(1-q_2)(1-P_3)$$

$$P(X=x_0, Y=y_1, Z=z_1) = r(1-q_2)P_3$$

$$P(X=x_1, Y=y_0, Z=z_1) = r q_2(1-P_4)$$

$$P(X=x_1, Y=y_1, Z=z_1) = r q_2 P_4$$

$$P(X, Y) = \sum_Z P(X|Z)P(Y|X, Z)$$

$$P(X=x_0, Y=y_0) = P(X=x_0|Z=z_0)P(Y=y_0|X=x_0, Z=z_0) + P(X=x_0|Z=z_1)P(Y=y_0|X=x_0, Z=z_1)$$

$$= (1-q_1)(1-P_1) + (1-q_2)(1-P_3)$$

$$P(X=x_0, Y=y_1) = (1-q_1)P_1 + (1-q_2)P_3$$

$$P(X=x_1, Y=y_0) = q_1(1-P_2) + q_2(1-P_4)$$

$$P(X=x_1, Y=y_1) = q_1P_2 + q_2P_4$$

$$P(X, Z) = P(Z)P(X|Z)$$

$$P(X=x_0, Z=z_0) = (1-r)(1-q_1) \quad P(X=x_1, Z=z_0) = (1-r)q_1$$

$$P(X=x_0, Z=z_1) = r(1-q_2) \quad P(X=x_1, Z=z_1) = r q_2$$

$$P(X=x_0, Z=z_1) = r(1-q_2) \quad P(X=x_1, Z=z_1) = r q_2$$

$$P(Y, Z) = \sum_X P(Z) P(Y | X, Z)$$

$$\begin{aligned} P(Y=y_0, Z=z_0) &= P(Z=z_0) P(Y=y_0 | X=x_0, Z=z_0) \\ &\quad + P(Z=z_0) P(Y=y_0 | X=x_1, Z=z_0) \\ &= (1-r)(1-p_1) + (1-r)(1-p_2) \end{aligned}$$

$$P(Y=y_0, Z=z_1) = r(1-p_3) + r(1-p_4)$$

$$P(Y=y_1, Z=z_0) = (1-r)p_1 + (1-r)p_2$$

$$P(Y=y_1, Z=z_1) = rp_3 + rp_4$$

E2  $P(y_1 | x_1) - P(y_1 | x_0)$

$$P(y|x) = \frac{P(x, y)}{P(x)} = \frac{P(x, y)}{P(z) P(x|z)}$$

①  $Z=z_1$

$$\begin{aligned} P(y|x) &= \frac{P(x|z=z_1) P(y|x, z=z_1)}{P(z=z_1) P(x|z=z_1)} \\ &= \frac{P(y|x, z=z_1)}{P(z=z_1)} \end{aligned}$$

$$P(y_1 | x_1) = \frac{P(y_1 | x_1, z_1)}{P(z_1)} = \frac{p_4}{r}$$

$$P(y_1 | x_0) = \frac{P(y_1 | x_0, z_1)}{P(z_1)} = \frac{p_3}{r}$$

$$P(y_1 | x_1) - P(y_1 | x_0) = \frac{p_4 - p_3}{r}$$

$$P(Y, z) = \sum_x P(z) P(Y | X, z)$$

$$\begin{aligned} P(Y=y_0, z=z_0) &= P(z=z_0) P(Y=y_0 | X=x_0, z=z_0) \\ &\quad + P(z=z_0) P(Y=y_0 | X=x_1, z=z_0) \\ &= (1-r)(1-p_1) + (1-r)(1-p_2) \end{aligned}$$

$$P(Y=y_0, z=z_1) = r(1-p_3) + r(1-p_4)$$

$$P(Y=y_1, z=z_0) = (1-r)p_1 + (1-r)p_2$$

$$P(Y=y_1, z=z_1) = rp_3 + rp_4$$

E2  $P(y_1 | x_1) - P(y_1 | x_0)$

$$P(y|x) = \frac{P(x, y)}{P(x)} = \frac{P(x, y)}{P(z) P(x|z)}$$

①  $z = z_1$

$$\begin{aligned} P(y|x) &= \frac{P(x|z=z_1) P(y|x, z=z_1)}{P(z=z_1) P(x|z=z_1)} \\ &= \frac{P(y|x, z=z_1)}{P(z=z_1)} \end{aligned}$$

$$P(y_1 | x_1) = \frac{P(y_1 | x_1, z_1)}{P(z_1)} = \frac{p_4}{r}$$

$$P(y_1 | x_0) = \frac{P(y_1 | x_0, z_1)}{P(z_1)} = \frac{p_3}{r}$$

$$P(y_1 | x_1) - P(y_1 | x_0) = \frac{p_4 - p_3}{r}$$

$$(2) \quad z = z_0$$

$$P(y|x) = \frac{P(x|z_0) P(y|x, z_0)}{P(z_0) P(x|z_0)} = \frac{P(y|x, z_0)}{P(z_0)}$$

$$P(y_1|x_1) = \frac{P_2}{1-r} \quad P(y_1|x_0) = \frac{P_1}{1-r}$$

$$P(y_1|x_1) - P(y_1|x_0) = \frac{P_2 - P_1}{1-r}$$

$$(3) \quad z = (z_0, z_1)$$

$$P(y|x) = \sum_z \frac{P(y|x, z) P(z)}{P(x)}$$

$$P(y_1|x_1) = \frac{P_2}{1-r} + \frac{P_4}{r}$$

$$P(y_1|x_0) = \frac{P_1}{1-r} + \frac{P_3}{r}$$

$$\begin{aligned} P(y_1|x_1) - P(y_1|x_0) &= \frac{P_2 - P_1}{1-r} + \frac{P_4 - P_3}{r} \\ &= \frac{(P_2 + P_3 - P_1 - P_4)r + P_4 - P_3}{(1-r)r} \end{aligned}$$

E3 ? not sure.

E4

$$P(y | do(x)) = \sum_z P(Y=y | X=x, Z=z) P(Z=z)$$

$$\begin{aligned} P(y_0 | do(x_0)) &= P(y_0 | x_0, z_0) P(z_0) + P(y_0 | x_0, z_1) P(z_1) \\ &= (1 - P_1)(1 - r) + (1 - P_3)r \end{aligned}$$

$$P(y_1 | do(x_0)) = P_1(1 - r) + P_3 r$$

$$\begin{aligned} P(y_0 | do(x_1)) &= P(y_0 | x_1, z_0) P(z_0) + P(y_0 | x_1, z_1) P(z_1) \\ &= (1 - P_2)(1 - r) + (1 - P_4)r \end{aligned}$$

$$P(y_1 | do(x_1)) = P_2(1 - r) + P_4 r$$

E5

$$ACE = P(y_1 | do(x_1)) - P(y_1 | do(x_0)) =$$

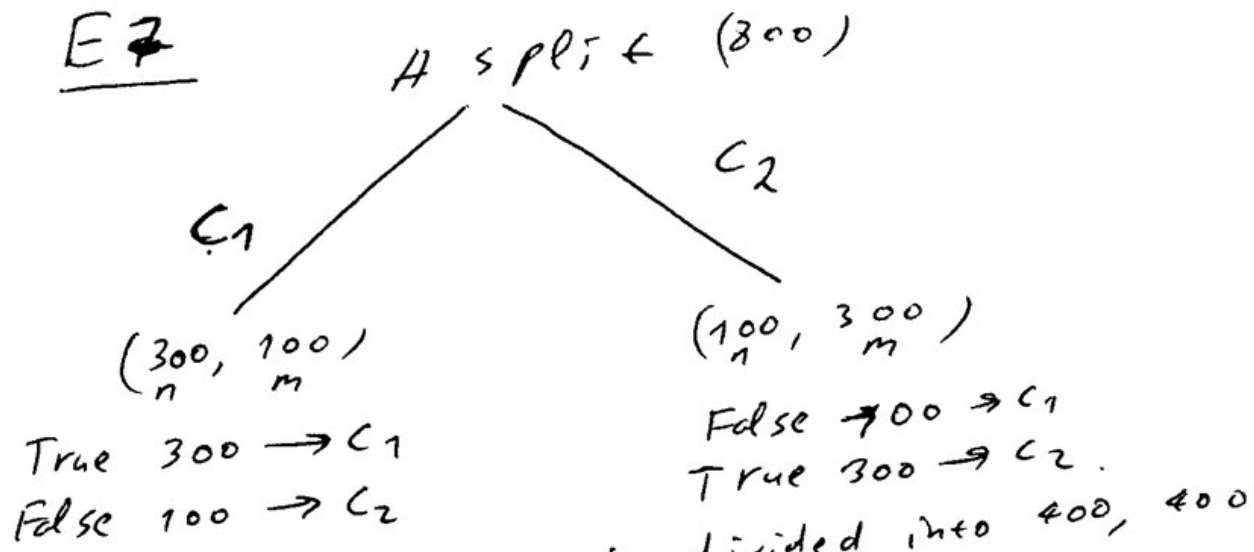
$$= P_2 - P_2 r + P_4 r - P_1 + P_1 r - P_3 r$$

$$= P_2 - P_1 + r(P_1 + P_4 - P_2 - P_3)$$

$$RD = \frac{(P_2 + P_3 - P_1 - P_4)r + P_4 - P_3}{(1 - r)r}$$

RD normalized with  $(1 - r)r$  and  $P_n$  values looks like negatively inversed.

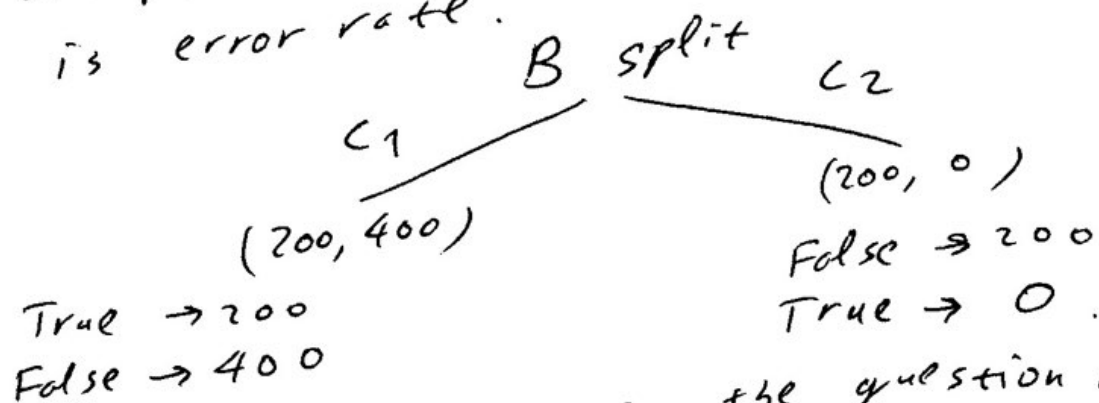
E2



All 800 data is divided into 400, 400  
Left branch is for predicting C1 and  
has 300 points from C1 and 100 points from C2  
Right branch is also same.

$$\text{So error rate} = \frac{FP + FN}{TP + TN + FP + FN} =$$
$$= \frac{100 + 100}{800} = \frac{1}{4} = 0,25$$

all False prediction divided by all values  
is error rate.



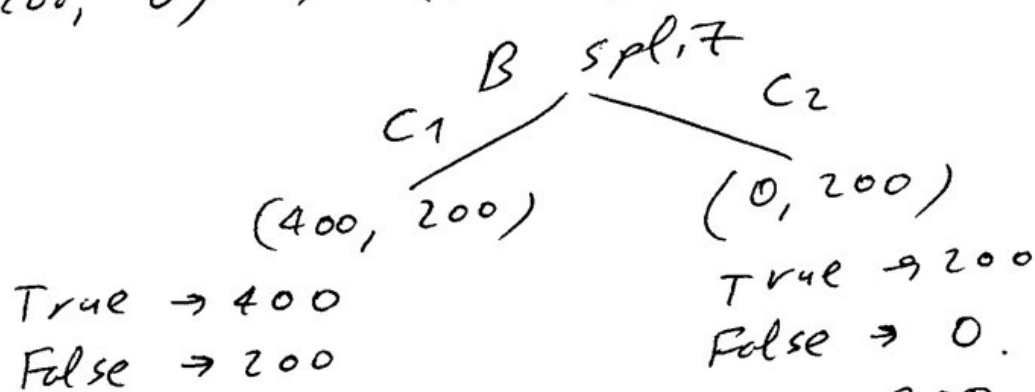
Here, I am assuming the question is  
wrong. There is no way error rate is equal to A.  
We have 600 points wrongly predicted  
in B

Where A predicted only 200 wrong.

I am going to flip the numbers

$$(200, 400) \rightarrow (400, 200)$$

$$(200, 0) \rightarrow (0, 200)$$



$$\text{error rate} = \frac{FP+FN}{TP+TN+FP+FN} = \frac{200}{800} = 0,25$$

so both error rate is same 0,25.

E 8. split A (using the table from lecture)

	C1	C2	
pred C1	300(n <sub>11</sub> )	100(n <sub>12</sub> )	400(n <sub>1+</sub> )
pred C2	100(n <sub>21</sub> )	300(n <sub>22</sub> )	400(n <sub>2+</sub> )
	400(n <sub>1+</sub> )	400(n <sub>2+</sub> )	n

$$\text{Gini}(\text{pred C1}) = 2 \cdot \frac{n_{11}}{n_{1+}} \cdot \frac{n_{12}}{n_{1+}} = 2 \cdot \frac{300 \cdot 100}{160000} = \frac{3}{8}$$

$$\text{Gini}(\text{pred C2}) = 2 \cdot \frac{n_{21}}{n_{2+}} \cdot \frac{n_{22}}{n_{2+}} = 2 \cdot \frac{100 \cdot 300}{160000} = \frac{3}{8}$$

$$\text{Gini} = \frac{400}{800} \cdot \frac{3}{8} + \frac{400}{800} \cdot \frac{3}{8} = \frac{3}{8}$$

split B

	C1	C2	
pred C1	400(n <sub>11</sub> )	200(n <sub>12</sub> )	600(n <sub>1+</sub> )
pred C2	0(n <sub>21</sub> )	200(n <sub>22</sub> )	200(n <sub>2+</sub> )
	400(n <sub>1+</sub> )	400(n <sub>2+</sub> )	n



$$\text{Gini}(\text{pred } c_1) = 2 \cdot \frac{n_{11}}{n_{+1}} \cdot \frac{n_{12}}{n_{+1}} = 2 \cdot \frac{400 \cdot 200}{360000} = \frac{4}{9}$$

$$\text{Gini}(\text{pred } c_2) = 2 \cdot \frac{n_{21}}{n_{+2}} \cdot \frac{n_{22}}{n_{+2}} = 0$$

$$\text{Gini} = \frac{n_{+1}}{n} \cdot \text{Gini}(\text{pred } c_1) = \frac{400}{800} \cdot \frac{4}{9} = \frac{2}{9}$$

$$\text{Gini}(A) = \frac{3}{8} > \text{Gini}(B) = \frac{2}{9}$$

since  $\text{Gini}(A)$  has more Gini, we are more uncertain about split A and more certain about split B.

$$\text{Entropy} = -\sum p_i \log p_i = -p_1 \log p_1 - p_2 \log p_2$$

$$\text{A split entropy}(\text{pred } c_1) = -\frac{n_{11}}{n_{+1}} \log \frac{n_{11}}{n_{+1}} - \frac{n_{12}}{n_{+1}} \log \frac{n_{12}}{n_{+1}}$$

$$= -\frac{300}{400} \log \frac{3}{4} - \frac{100}{400} \log \frac{1}{4} = 0,31127 + 0,50000 = 0,81127$$

$$\text{entropy}(\text{pred } c_2) = -\frac{n_{21}}{n_{+2}} \log \frac{n_{21}}{n_{+2}} - \frac{n_{22}}{n_{+2}} \log \frac{n_{22}}{n_{+2}}$$

$$= -\frac{100}{400} \log \frac{1}{4} - \frac{300}{400} \log \frac{3}{4} = 0,81127$$

$$\text{entropy} = \frac{n_{+1}}{n} \text{entropy}(\text{pred } c_1) + \frac{n_{+2}}{n} \text{entropy}(c_2)$$

$$= \frac{1}{2} (0,81127 + 0,81127) = 0,81127$$

$$\text{B split entropy}(\text{pred } c_1) = -\frac{4}{6} \log \frac{4}{6} - \frac{2}{6} \log \frac{2}{6} =$$

$$= 0,38997 + 0,52832 = 0,91829$$

$$\text{entropy}(\text{pred } c_2) = 0$$

$$\text{entropy} = \frac{1}{2} \cdot 0,91829 = 0,45914$$

$$\text{entropy } A = 0,81127 > \text{entropy } B = 0,45914$$

E9  $f_s(x_s) = E_{x_c}[f(x_s, x_c)]$

$$f(x) = f_1(x_s) + f_2(x_c)$$

$$f_s(x_s) = E_{x_c}[f(x_s, x_c)] = E_{x_c}[f_1(x_s, x_c) + f_2(x_c, x_c)]$$

$$= \underbrace{E_{x_c}[f_1(x_s + x_c)]}_{f_1(x_s)} + E_{x_c}[f_2(x_c, x_c)]$$

$$= f_1(x_s) + \underbrace{E_{x_c}[f_2(x_c, x_c)]}_{\text{doesn't depend on } x_s \text{ so we can assume it is constant}}$$

$$= f_1(x_s) + C$$

E10  $f(x) = f_1(x_s) \cdot f_2(x_c)$

$$f_s(x_s) = E_{x_c}[f(x_s, x_c)] = E_{x_c}[f_1(x_s, x_c) f_2(x_c, x_c)]$$

$$= \underbrace{E_{x_c}[f_1(x_s, x_c)]}_{f_1(x_s)} E_{x_c}[f_2(x_c, x_c)]$$

$$= f_1(x_s) \underbrace{E_{x_c}[f_2(x_c, x_c)]}_{\text{doesn't depend on } x_s \text{ so assume constant}}$$

$$= f_1(x_s) \cdot C$$