

Nama : M. Faishal Muntaz

NIM : 664180093

1. a. - Tidak lengkap (missing Value) pada kolom Height
- Tidak konsisten pada satuan kolom weight (kg, lbs)
  - Tidak konsisten pada kolom Covid-19 Result
  - Tidak konsisten pada kolom Blood group
  - Terdapat noisy pada age yang berbeda

Cara mengatasinya :

- Menyamakan satuan pada kolom yang tidak konsisten
- Pada missing Value dimasukan nilai rata-rata pada kolom tersebut
- Pada noisy melakukan pembersihan data.

b. Equal - width (weight)

$$w = (\max - \min) / N$$

=> membagi 3 kategori

$$(95 - 45) / 3 = 13,3$$

$$\text{Range 1 (rendah)} : 45 + 13,3 = 58,3 \quad (0 - 58,3)$$

$$\text{Range 2 (Sedang)} : 58,3 + 13,3 = 71,6 \quad (58,3 - 71,6)$$

$$\text{Range 3 (Tinggi)} : 71,6 + 13,3 = 84,9 \approx 85 \quad (71,6 - 95)$$

## 2. Single linkage Clustering

	$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$	$T_9$	$T_{10}$
$J_1$	0	1	0	0	0	0	1	1	0	0
$J_2$	1	0	1	0	1	1	0	1	0	1
$J_3$	1	1	0	0	0	1	1	0	0	1

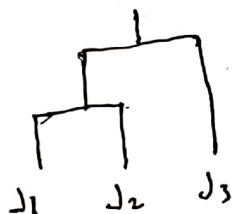
\* Sim Jaccard:  $\frac{a}{a+b+c}$

Sim  $J_1, J_2 = \frac{1}{1+2+5} = \frac{1}{8}$  - karena terkecil maka dipilih

Sim  $J_1, J_3 = \frac{2}{2+1+3} = \frac{2}{6}$

Sim  $J_2, J_3 = \frac{3}{5+3+2} = \frac{3}{10}$

Dendogramnya  $\Rightarrow$



## 3. $a(0,1)$ , $b(1,1)$ , $c(1,2)$

\* Manhattan =  $|x_1 - x_0| + |y_1 - y_0|$

	a	b	c
a	-	1	2
b		-	1
c			-

\* Dist. Jarak terdekat ke-2

dist(a) = c) = 2

dist(b) = a) = 1

dist(c) = a) = 2

\* titik dgn Jarak  $\leq$  dist

$N_2(a) = \{b, c\} = 2$

$N_2(b) = \{a, c\} = 2$

$N_2(c) = \{a, b\} = 2$

\*  $lrd(a)$

$lrd(a) = \frac{2}{1+2} = \frac{2}{3} = 0,66$

$lrd(b) = \frac{2}{2+2} = \frac{2}{4} = 0,5$

$lrd(c) = \frac{2}{2+1} = \frac{2}{3} = 0,66$

\*  $LoF(a)$

$LoF(a) = (0,5 + 0,66) * (1+2) = 3,48$

$LoF(b) = (0,66 + 0,66) * (2+2) = 5,28$

$LoF(c) = (0,66 + 0,5) * (2+1) = 3,48$

o o Top 1 outlier adalah b dengan

$LoF(b) = 5,28$

4. \* Temperatur Badan  
\* Normal (36-37)

Temperatur - Normal	
Ya	2
Tidak	2

$$\text{Gini} = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2$$

$$= 0,5$$

Temperatur - Tidak normal	
Ya	3
Tidak	1

$$\text{Gini} = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2$$

$$= 0,375$$

- Gini Split

$$= \frac{4}{6} * 0,5 + \frac{4}{6} * 0,375$$

$$= 0,4375$$

\* Saturasi oksigen  
rendah jika  $\leq 90\%$

Saturasi - rendah	
Ya	5
Tidak	0

$$\text{Gini} = 1 - \left(\frac{5}{5}\right)^2 - (0)^2$$

$$= 0$$

Saturasi - tinggi	
Ya	0
Tidak	3

$$\text{Gini} = 1 - (0)^2 - \left(\frac{3}{3}\right)^2$$

$$= 0$$

Gini Split

$$= \frac{5}{6} * 0 + \frac{3}{6} * 0$$

$$= 0$$

\* Rapid Test - reaktif

Rapid Test - reaktif	
Ya	3
Tidak	1

$$\text{Gini} = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2$$

$$= 0,375$$

Rapid Test - non reaktif

Rapid Test - non reaktif	
Ya	2
Tidak	2

$$\text{Gini} = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2$$

$$= 0,5$$

Gini Split

$$= \frac{4}{6} * 0,375 + \frac{4}{6} * 0,5$$

$$= 0,4375$$

Best Splitnya adalah Saturasi oksigen karena ~~adalah~~ Gini Split yang paling kecil, yaitu 0

5. a

		Predicted		
		+	-	
Actual	+	250	250	500
	-	50	1450	1500
Total		300	1700	2000

b.  $TP = 250$

$TN = 1450$

$FP = 50$

$FN = 250$

c.  $Accuracy = \frac{250 + 1450}{2000} = 0,85$

$Precision = \frac{250}{250 + 50} = 0,83$

$Recall = \frac{250}{250 + 250} = 0,5$

$F-measure = \frac{2 \times 0,83 \times 0,5}{0,83 + 0,5} = 0,624$

d. Model kurang baik dikarenakan nilai recall dan F-measure kecil walaupun tingkat accuracynya tinggi