

Question #1

A) employ IR to answer the following
which attribute produces lowest error rate for predicting
the outcome of >10 chinups?

First, build the table:

Weight	Waist		HeartRate		Chinups		yes	no
	yes	no	yes	no	fast	slow		
heavy	3/11	6/9	thin	9/11	2/9	fast	6/11	4/9
light	8/11	3/9	wide	2/11	7/9	slow	5/11	5/9

from the table, clearly waist has lowest error rate

rule #1: if waist=thin then >10 chinups = yes ($\frac{9}{11}$)

rule #2: if waist=wide then >10 chinups = no ($\frac{7}{9}$)

B) what is the error rate of rules from (A)?

$$\text{rule 1: } \frac{9}{11} = 0.818, 1 - 0.818 = 0.182 \text{ or } 18\%$$

$$\text{rule 2: } \frac{7}{9} = 0.778, 1 - 0.778 = 0.222 \text{ or } 22\%$$

Question #2 Name Bayes

A) What is probability >10 chinups = yes given weight = heavy, waist = wide, and heartrate = slow?

$$\text{Bayes Theorem: } p(X|Y) = \frac{p(Y|X)p(X)}{p(Y)}$$

We want to know $p(>10 \text{ chinups} = \text{yes} | \text{heavy, wide, slow})$

$$\Rightarrow p(\text{yes} | \text{heavy, wide, slow}) = \frac{p(\text{heavy, wide, slow} | \text{yes})p(\text{yes})}{p(\text{heavy, wide, slow})}$$

using marginal probability:

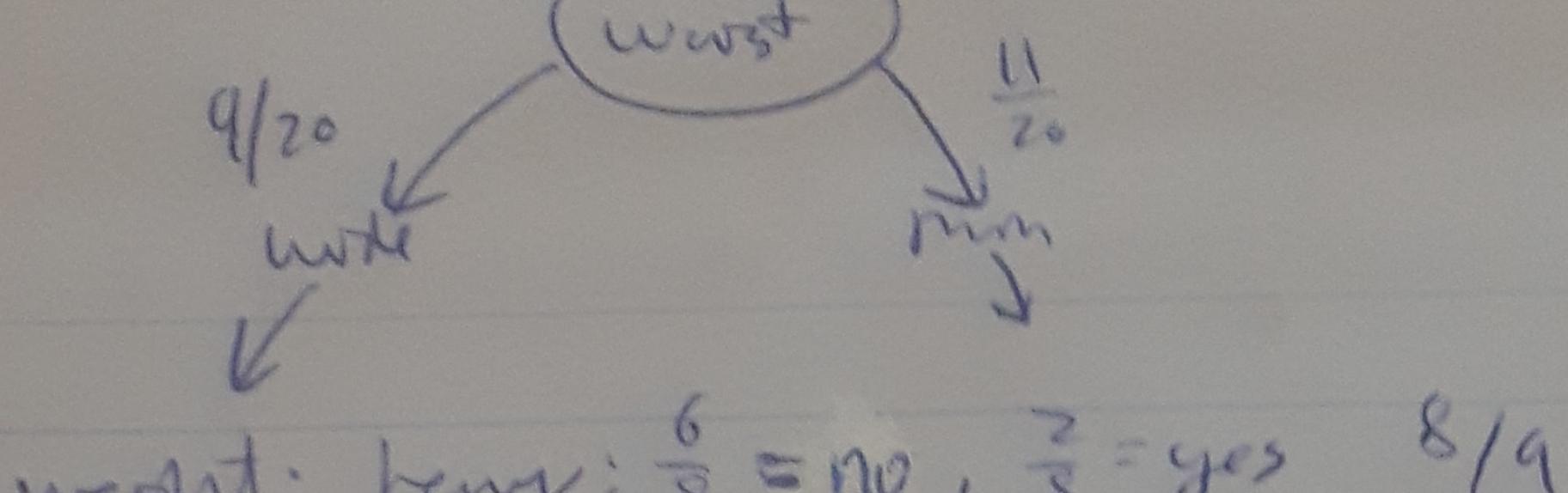
$$p(\text{heavy, wide, slow}) = p(\text{heavy, wide, slow} | \text{yes})p(\text{yes}) + p(\text{heavy, wide, slow} | \text{no})p(\text{no})$$

$$\Rightarrow p(\text{yes} | \text{heavy, wide, slow}) = \frac{p(\text{heavy, wide, slow} | \text{yes})p(\text{yes})}{p(\text{heavy, wide, slow} | \text{yes})p(\text{yes}) + p(\text{heavy, wide, slow} | \text{no})p(\text{no})}$$

assuming independence, $p(\text{heavy, wide, slow} | \text{yes}) = p(\text{heavy} | \text{yes})p(\text{wide} | \text{yes})p(\text{slow} | \text{yes})$

$$\Rightarrow p(\text{yes} | \text{heavy, wide, slow}) = \frac{\frac{3}{11} \cdot \frac{2}{11} \cdot \frac{11}{20}}{\left(\frac{3}{11} \cdot \frac{2}{11} \cdot \frac{11}{20}\right) + \left(\frac{6}{11} \cdot \frac{7}{9} \cdot \frac{9}{20}\right)} = 0.0373$$

B) which values of weight, waist, and heartrate yield highest probability of >10 chinups = yes? answer: simply find the highest fractions for each attribute: light, thin, fast



$$\text{weight: heavy: } \frac{6}{11} = \text{no}, \frac{2}{11} = \text{yes} \quad \frac{8}{9}$$

$$\text{light: } \frac{1}{11} = \text{no}, \frac{9}{9} = \text{yes} \quad \frac{1}{9}$$

$$\text{heartrate: slow: } \frac{5}{11} = \text{no}, \frac{6}{9} = \text{yes}, \frac{5}{9}$$

$$\text{fast: } \frac{2}{11} = \text{no}, \frac{7}{9} = \text{yes}, \frac{4}{9}$$

$$Gini(\text{weight}) = 1 - \left(\frac{8}{9} \left(\frac{6}{11} \right)^2 + \left(\frac{2}{11} \right)^2 \right) + \left(\frac{1}{9} \left(\frac{1}{1} \right)^2 + \left(\frac{9}{9} \right)^2 \right) = 0.338$$

$$Gini(\text{heartrate}) = 0.222 \Rightarrow \text{heartrate}$$

$$\text{weight: heavy: } \frac{1}{11}, \frac{1}{1} = \text{yes}, \frac{0}{1} = \text{no}$$

$$\text{light: } \frac{10}{11}, \frac{8}{10} = \text{yes}, \frac{2}{10} = \text{no}$$

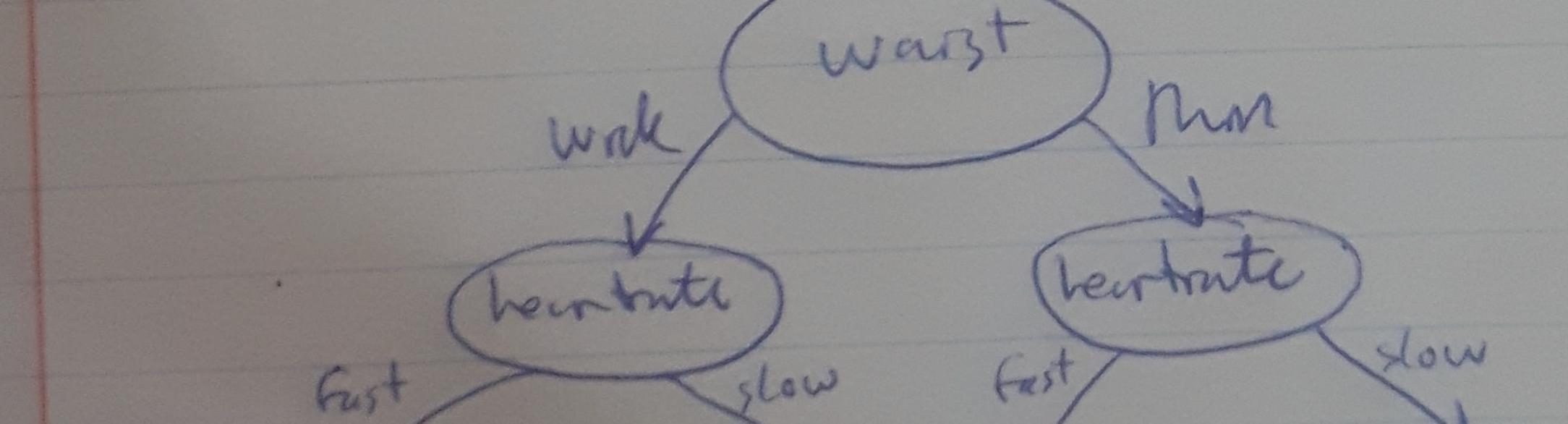
$$Gini(\text{weight}) = 0.2909$$

⇒ heartrate

$$\text{heartrate: fast: } \frac{6}{11}, \frac{4}{6} = \text{yes}, \frac{2}{6} = \text{no}$$

$$\text{slow: } \frac{5}{11}, \frac{5}{5} = \text{yes}, \frac{0}{5} = \text{no}$$

$$Gini(\text{heartrate}) = 0.242$$



$$\text{weight: heavy: } \frac{1}{11}, \frac{1}{1} = \text{yes}, \frac{0}{1} = \text{no}$$

$$\text{light: } \frac{10}{11}, \frac{8}{10} = \text{yes}, \frac{2}{10} = \text{no}$$

$$Gini(\text{weight}) = 0.2909$$

⇒ heartrate

$$\text{weight: heavy: } \frac{1}{11}, \frac{1}{1} = \text{yes}, \frac{0}{1} = \text{no}$$

$$\text{light: } \frac{10}{11}, \frac{8}{10} = \text{yes}, \frac{2}{10} = \text{no}$$

$$Gini(\text{weight}) = 0.2909$$

⇒ heartrate

#3 Use divide and conquer to construct a decision tree on the dataset, using Gini index as splitting criterion, to answer the following:

A) which attribute the best initial split? $Gini = 1 - \sum_i P_i^2$

$$\text{weight: heavy: } \frac{3}{9} = \text{yes}, \frac{6}{9} = \text{no}$$

$$\text{light: } \frac{8}{11} = \text{yes}, \frac{3}{11} = \text{no}$$

$$Gini(\text{weight}) = 1 - \left(\frac{9}{20} \left(\frac{3}{9} \right)^2 + \left(\frac{6}{9} \right)^2 \right) + \left(\frac{11}{20} \left(\frac{8}{11} \right)^2 + \left(\frac{3}{11} \right)^2 \right) = 0.4182$$

$$\text{waist: wide: } \frac{2}{9} = \text{yes}, \frac{7}{9} = \text{no}$$

$$\text{thin: } \frac{9}{11} = \text{yes}, \frac{2}{11} = \text{no}$$

$$Gini(\text{waist}) = 0.3192$$

$$\text{heartrate: fast: } \frac{6}{10} = \text{yes}, \frac{4}{10} = \text{no}$$

$$\text{slow: } \frac{5}{10} = \text{yes}, \frac{5}{10} = \text{no}$$

$$Gini(\text{heartrate}) = 0.49$$

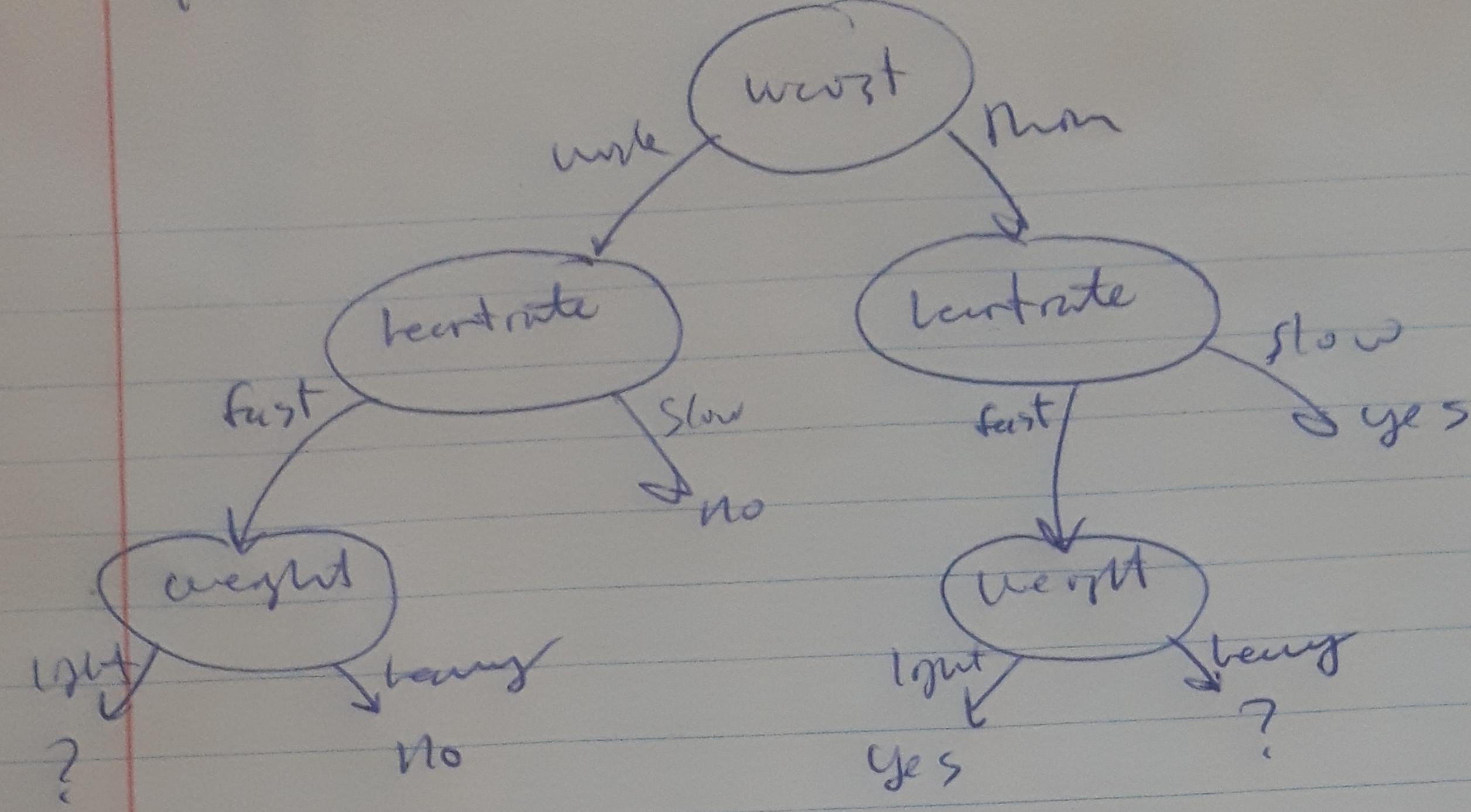
split on waist as root.

Waist		Weight		HR		>10	
thin	wide	light	heavy	slow	fast	no	yes
light	fast	yes	no	no	no	no	no
heavy	slow	yes	no	no	no	no	no
light	fast	no	yes	no	no	yes	yes
heavy	fast	no	yes	no	no	yes	yes
light	slow	no	yes	no	no	yes	yes
heavy	slow	no	yes	no	no	yes	yes
light	fast	yes	no	no	no	yes	yes
heavy	fast	yes	no	no	no	yes	yes
light	slow	yes	no	no	no	yes	yes
heavy	slow	yes	no	no	no	yes	yes

because we split,
we create 2 sub-sets of data based
on instances that
had waist=wide
or waist=thin

We then calculate
gini on these
sub-sets to keep
building tree

Final tree



Rules:

- A → if waist = wide and heartrate = slow then >10 = no
- B → if waist = thin and heartrate = slow then >10 = yes
- C → if waist = wide and heartrate = fast and weight = heavy then >10 = no
- D → if waist = thin and heartrate = fast and weight = light then >10 = yes

B) compare to IR

rule #1 from IR is a more general version of rule B from the tree

rule #2 from IR is a more general version of rule A from the tree