

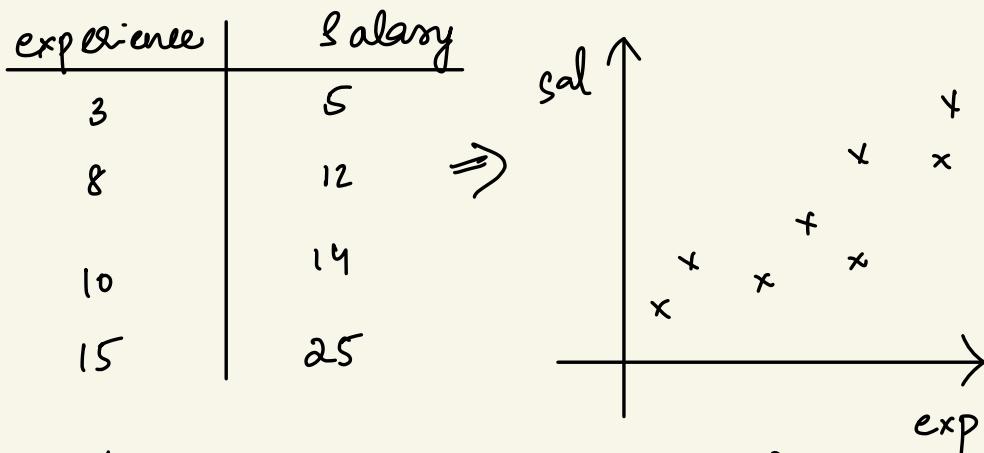
Decision Trees

About Parametric and non-parametric algorithms

Parametric Algorithms :-

- makes assumptions
- fixed parameters
- requires less data to train

e.g.: consider the following data:



Linear data $\Rightarrow y = f(x)$
 $\text{salary} = f(\text{experience})$

If done using linear regression:

$$y = mx + b \rightarrow \text{assuming it's a line}$$

\downarrow \downarrow
 fixed parameters

Non-parametric algorithms:-

- does not assume
- needs more data to train
- chances of over fitting
 - gives importance to all choices
 - might consider noise also.

e.g.: Decision Tree

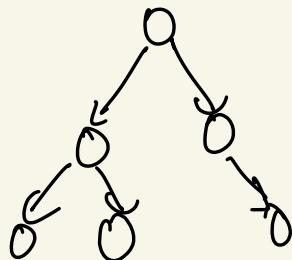
About Classification

- dividing dataset into categories
- in supervised → labels given.

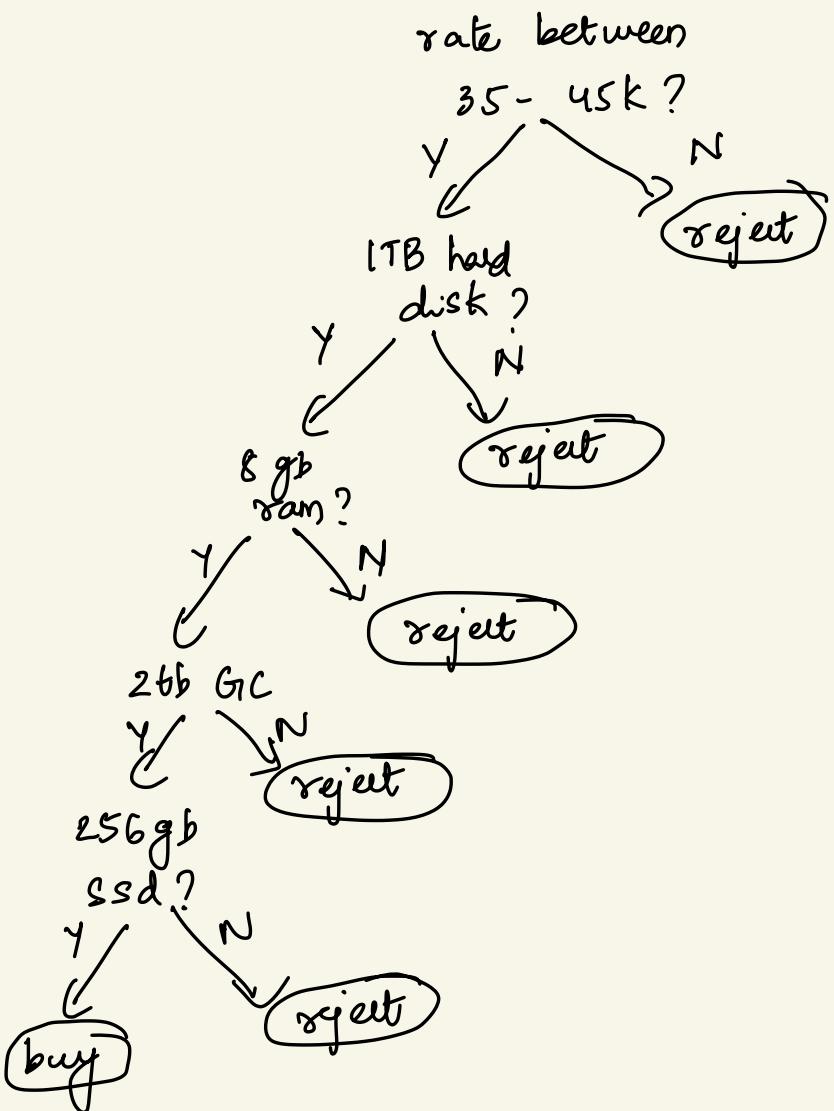
Decision Trees Algorithm

Transforming data into a tree structure.

- root node
 - no incoming
 - outgoing
- internal node
 - incoming & outgoing
- leaf node - incoming, no outgoing, class label predicted here.



Give example of laptop buying :-

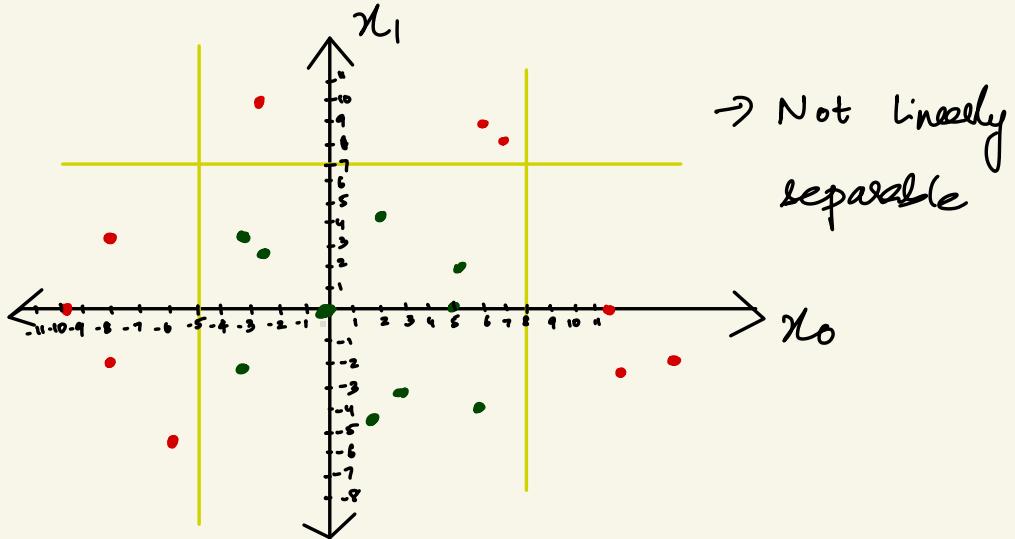


- it is hierarchical
- recursive splits used
- each decision node implements a function $f_m(x)$ with discrete outcomes labelling the branches
- process starts at root
- repeated recursively until leaf node is hit
- value at leaf node is the output
- tree structure is not fixed
- tree grows during learning depending on complexity

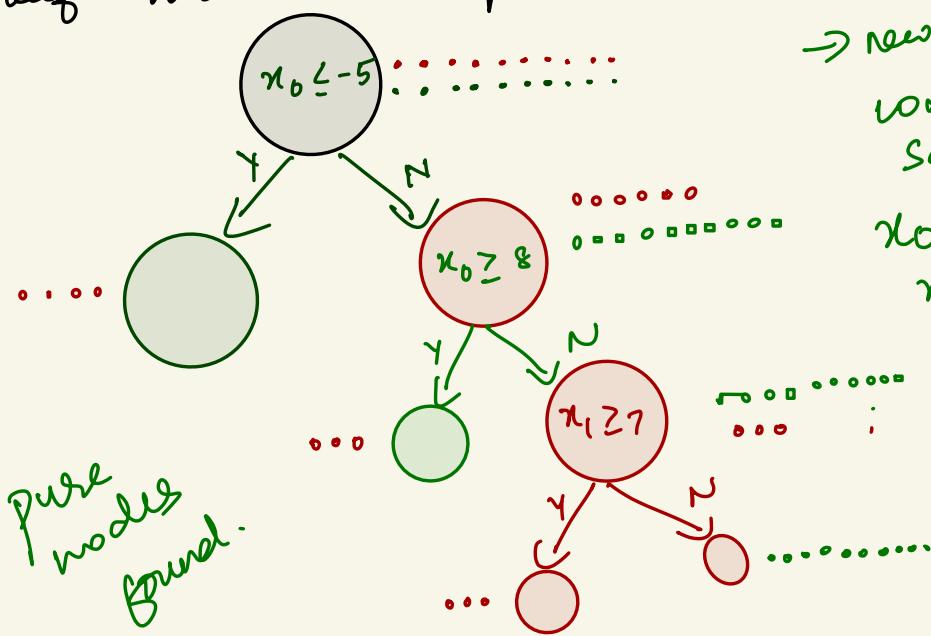
→

Let's consider a random dataset:

→ contains 2 features $\rightarrow x_0, x_1$



→ recursively splits a dataset to obtain leaf nodes that's pure.



What if we do not get pure nodes?

Impurity is calculated using entropy.

→ Entropy measures impurity of collection of samples

→ depends on the proportion of classes

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

So if we have 2 classes \Rightarrow positive & negative:

$$E(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

\downarrow
proportion of
the samples

\downarrow
proportion
of -ve
samples

proportion means no of +ve samples
out of total samples

Examples:

$$\begin{aligned} E(14+, 0-) &= -\frac{14}{14} \log_2\left(\frac{14}{14}\right) - 0 \log_2(0) \\ &= -1 \log_2(1) - 0 \log_2(0) \\ &= -1 \times 0 - 0 \Rightarrow \underline{\underline{0}} \end{aligned}$$

$$\begin{aligned} E(7+, 7-) &= -\frac{7}{14} \log_2\left(\frac{7}{14}\right) - \frac{7}{14} \log_2\left(\frac{7}{14}\right) \\ &= -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \\ &= -\frac{1}{2} \times -1 - \frac{1}{2} \times -1 \\ &= \frac{1}{2} + \frac{1}{2} = \underline{\underline{1}} \end{aligned}$$

Example dataset :-

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

If we want to find entropy of this dataset :

$$\text{No. of +ve} = 9$$

$$\text{No. of -ve} = 5$$

$$\text{Total rows} = 14$$

So with eqn,

$$E(9+, 5-) = - \frac{9}{14} \log_2 \left(\frac{9}{14} \right) -$$

$$\frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

$$\approx \underline{\underline{0.9}} \Rightarrow$$

- impurity measure is done between 0 & 1
- 0.9 → highly impure.

Example - 1

Decision Tree algorithms - Entropy & Information Gain

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Step 1: Decide root node.

Method :-

- find information gain of all attributes
- attribute with highest IG will become root node.
- Calculate impurity

Impurity measure → using Entropy.

$$E(S) = \sum_{i=1}^c - P_i \log_2 P_i$$

Here → $P_+ \log_2 P_+ - P_- \log_2 P_-$

step 1: Entropy of entire dataset

$$S\{9+, 5-\} =$$

$$-\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$\Rightarrow 0.94$$

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Attribute : Outlook

Possible values - Sunny,

Overcast, rain

$$S_{\text{Sunny}} = [2+, 3-]$$

$$= - \frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \Rightarrow 0.97$$

$$S_{\text{Overcast}} = [4+, 0-] \Rightarrow 0$$

$$S_{\text{rain}} = [3+, 2-] \Rightarrow - \frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \Rightarrow 0.97$$

After entropy calculation \rightarrow Calculate $|G|$

$$G_{\text{air}}(S, \text{outlook}) =$$

$$\text{Entropy}(S) - \sum_{v \in \{\text{sunny, rain, overcast}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\Rightarrow \text{Entropy}(S) - \frac{5}{14} \text{Entropy}(S_{\text{sunny}})$$

$$- \frac{4}{14} \text{Entropy}(S_{\text{overcast}})$$

$$- \frac{5}{14} \text{Entropy}(S_{\text{rain}})$$

$$\Rightarrow 0.94 - \frac{5}{14} \times 0.97 - \frac{4}{14} \times 0 = -\frac{5}{14} \times 0.97$$
$$\Rightarrow \underline{\underline{0.2464}}$$

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Attribute :- Temp

Possible values -
Hot, Mild, Cool

$$S_{\text{Hot}} = [2+, 2-]$$

$$\Rightarrow \underline{\underline{1}}$$

$$S_{\text{Mild}} = [4+, 2-]$$

$$\Rightarrow -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = \underline{\underline{0.91}}$$

$$S_{\text{cool}} = [3+, 1-] \Rightarrow -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = \underline{\underline{0.81}}$$

$$\text{Gain}(S, \text{Temp}) = \text{Entropy}(S) - \sum_{\text{values}} \frac{|S_{\text{value}}|}{|S|} \text{Entropy}(S_{\text{value}})$$

$$= \text{Entropy}(S) - \frac{4}{14} E(S_{\text{hot}}) - \frac{6}{14} E(S_{\text{mild}}) - \frac{4}{14} E(S_{\text{cool}})$$

$$= 0.94 - \frac{4}{14} \cdot 1.0 - \frac{6}{14} \cdot 0.91 - \frac{4}{14} \cdot 0.81 = \underline{\underline{0.028}}$$

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Attribute: Humidity

Possible values = high, normal

$$S_{\text{high}} = [3+, 4-]$$

$$E(S_{\text{high}})$$

$$= -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7}$$

$$= 0.98$$

$$S_{\text{normal}} = [6+, 1-]$$

$$E(S_{\text{normal}}) = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0.59$$

$$\text{Gain}(S, \text{humidity}) = E(S) - \sum_{\text{outcomes}} \frac{|S_{\text{out}}|}{|S|} E(S_{\text{out}})$$

$$= E(S) - \frac{7}{14} E(S_{\text{high}}) - \frac{7}{14} E(S_{\text{normal}})$$

$$= 0.94 - \frac{7}{14} \cdot 0.98 - \frac{7}{14} \cdot 0.59 = 0.15$$

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Attribute : Wind

Possible values - weak, strong

$$S_{\text{strong}} = \{3+, 3-\}$$

$$E(S_{\text{strong}}) = \underline{\underline{1}}$$

$$S_{\text{weak}} = [6+, 2-]$$

$$E(S_{\text{weak}}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = \underline{\underline{0.81}}$$

$$\text{Gain}(S, \text{wind}) = E(S) - \sum_{(\text{weak}, \text{strong})} \frac{|S_{\text{o}}|}{|S|} E(S_{\text{o}})$$

$$= E(S) - \frac{6}{14} E(S_{\text{strong}}) - \frac{8}{14} E(S_{\text{weak}})$$

$$= 0.94 - \frac{6}{14} \times 1 - \frac{8}{14} \times 0.81 = \underline{\underline{0.047}}$$

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Comparing IGs:

$$\text{Gain}(S, \text{outlook}) = 0.2464$$

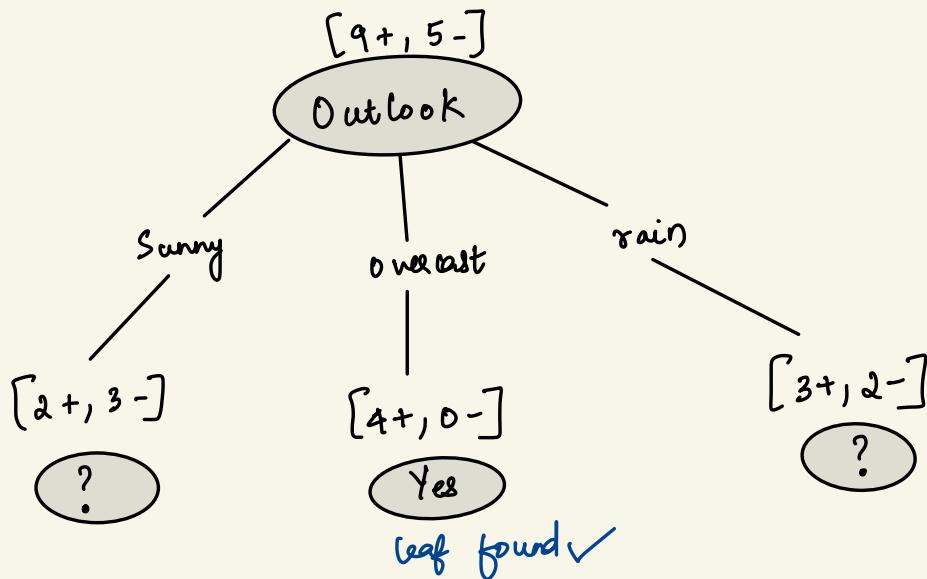
$$\text{Gain}(S, \text{temp}) = 0.028$$

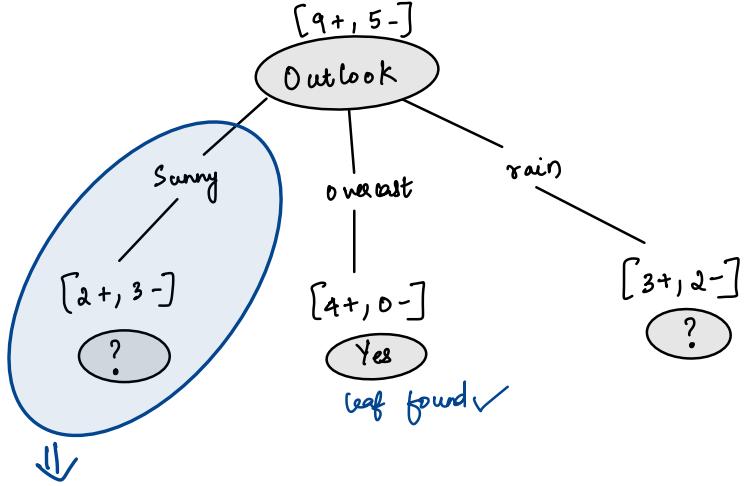
$$\text{Gain}(S, \text{humid}) = 0.15$$

$$\text{Gain}(S, \text{wind}) = 0.047$$

Humidity

high normal





take dataset with sunny values

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



Day	Temp	Humidity	Wind	Play
D1	Hot	High	weak	No
D2	Hot	High	Strong	No
D8	Mild	High	weak	No
D9	Cool	normal	weak	Yes
D11	Mild	normal	Strong	Yes

Day	Temp	Humidity	Wind	Play
D1	Hot	High	weak	No
D2	Hot	High	Strong	No
D3	Mild	High	weak	No
D4	Cool	Normal	weak	Yes
D5	Mild	Normal	Strong	Yes

Repeat steps here again:

Consider this as whole dataset.

$$\begin{aligned}
 E(S) &= E(S_{\text{Sunny}}) = [2+, 3-] \\
 &= -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \\
 &= \underline{\underline{0.97}}
 \end{aligned}$$

Consider attributes:

Attribute :- Temp \rightarrow Possible values - hot, mild, cool

$$\left. \begin{array}{l}
 S_{\text{hot}} \rightarrow [0+, 2-] \Rightarrow E(S_{\text{hot}}) = \underline{\underline{0}} \\
 S_{\text{mild}} \rightarrow [1+, 1-] \Rightarrow E(S_{\text{mild}}) = \underline{\underline{1}} \\
 S_{\text{cool}} \rightarrow [1+, 0-] \Rightarrow E(S_{\text{cool}}) = \underline{\underline{0}}
 \end{array} \right\} \text{gain} = \underline{\underline{0.57}}$$

Day	Temp	Humidity	wind	Play
D1	Hot	High	weak	No
D2	Hot	High	Strong	No
D8	Mild	High	weak	No
D9	Cool	normal	weak	Yes
D11	Mild	normal	strong	Yes

Attribute - Humidity

Values - High, normal

$$S_{\text{high}} = [0+, 3-] \Rightarrow E(S_{\text{high}}) = \underline{\underline{0}}$$

$$S_{\text{normal}} = [2+, 0-] \Rightarrow E(S_{\text{normal}}) = \underline{\underline{0}}$$

$$\text{Gain}(S_{\text{sunny}}, \text{humidity}) = E(S) - \sum_{v \in \{ \text{high, normal} \}} \frac{|S_v|}{|S|} E(S_v)$$

$$\text{Gain} = 0.97 - 0 \Rightarrow \underline{\underline{0.97}}$$

Day	Temp	Humidity	Wind	Play
D1	Hot	High	weak	No
D2	Hot	High	Strong	No
D8	Mild	High	weak	No
D9	Cool	normal	weak	Yes
D11	Mild	normal	strong	Yes

Attribute = Wind

values \rightarrow Strong, weak

$$S_{\text{Strong}} = [1+, 1-] = E(S_{\text{Strong}}) = 1$$

$$S_{\text{weak}} = [1+, 2-] = E(S_{\text{weak}})$$

$$\Rightarrow -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.91$$

$$\text{gain} = 0.97 - \frac{2}{5} \cdot 1 - \frac{3}{5} \cdot 0.91 = 0.019$$

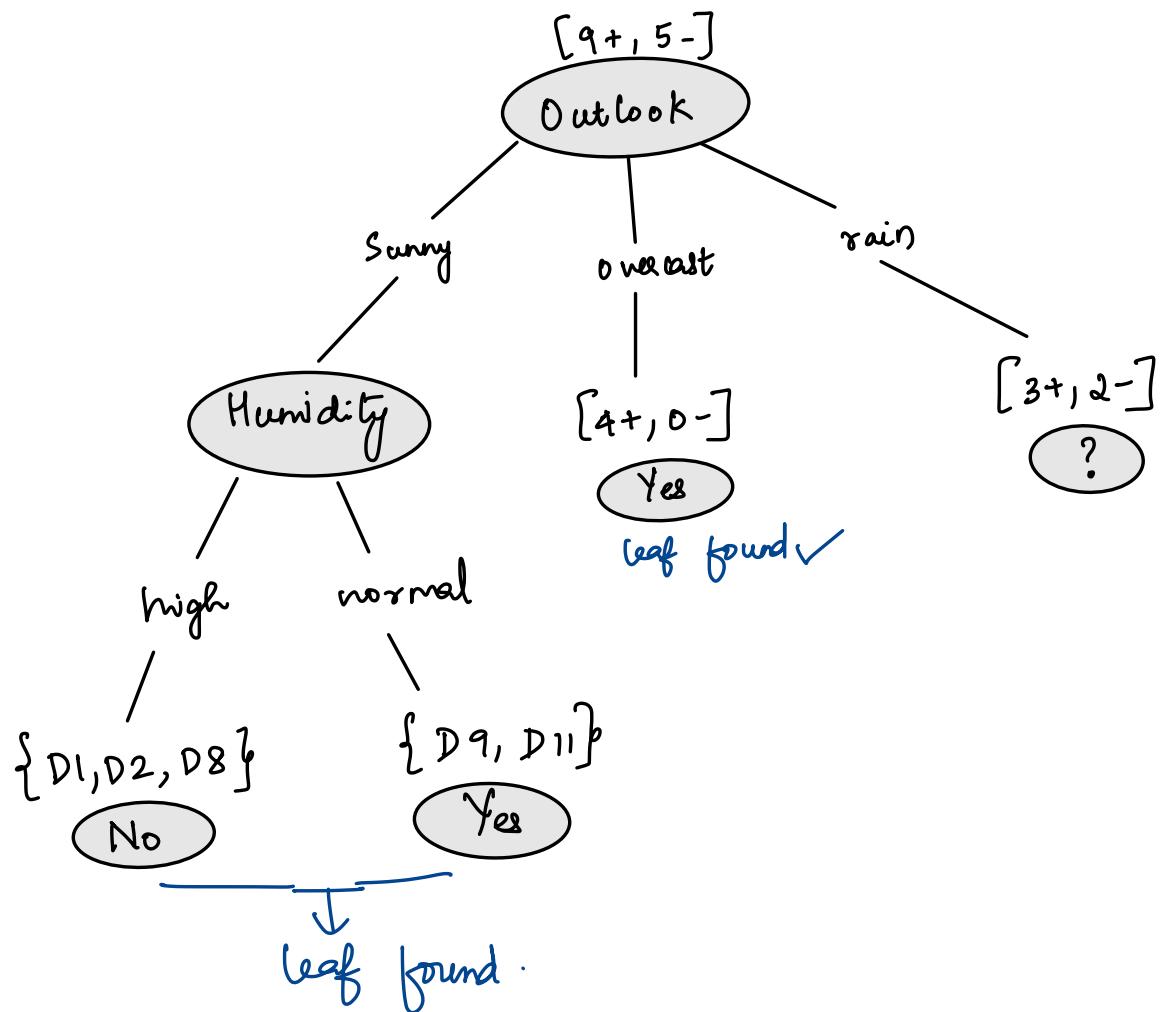
Day	Temp	Humidity	wind	Play
D1	Hot	High	weak	No
D2	Hot	High	Strong	No
D8	Mild	High	weak	No
D9	Cool	normal	weak	Yes
D11	Mild	normal	strong	Yes

Compare gains:-

$$\text{gain}(S_{\text{sunny}}, \text{temp}) = 0.57$$

$$\text{gain}(S_{\text{sunny}}, \text{humid}) = 0.97$$

$$\text{gain}(S_{\text{sunny}}, \text{wind}) = 0.019$$



Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

D4, D5, D6, D10, D14
 rain rows

Day	Temp	Humidity	Wind	Play
D4	Mild	high	weak	Yes
D5	Cool	normal	weak	Yes
D6	Cool	normal	Strong	No
D10	mild	normal	weak	Yes
D14	mild	high	Strong	No

Day	Temp	Humidity	Wind	Play
D4	Mild	high	weak	Yes
D5	Cool	normal	weak	Yes
D6	Cool	normal	Strong	No
D10	mild	normal	weak	Yes
D14	mild	high	strong	No

Repeat steps :-

Consider this as dataset :-

$$S_{\text{rain}} \Rightarrow [3+, 2-]$$

$$\begin{aligned} E(S_{\text{rain}}) &= -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \\ &= \underline{\underline{0.97}} \end{aligned}$$

Attributes : Temp

$$S_{\text{hot}} = [0+, 0-] \Rightarrow E(S_{\text{hot}}) = \underline{\underline{0}} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{gain} = 0.019$$

$$S_{\text{mild}} = [2+, 1-] \Rightarrow E(S_{\text{mild}}) = \underline{\underline{0.91}}$$

$$S_{\text{cool}} = [1+, 1-] \Rightarrow E(S_{\text{cool}}) = \underline{\underline{1}}$$

Day	Temp	Humidity	Wind	Play
D4	Mild	high	weak	Yes
DS	Cool	normal	weak	Yes
D6	Cool	normal	Strong	No
D10	mild	normal	weak	Yes
D14	mild	high	Strong	No

Attribute : Humidity

$$S_{\text{high}} = \{1+, 1-\} = E(S_{\text{high}}) = \underline{\underline{1}}$$

$$S_{\text{normal}} = \{2+, 1-\} = E(S_{\text{normal}}) = \underline{\underline{0.91}}$$

$$\text{gain}(S_{\text{rain}}, \text{humid}) = \underline{\underline{0.019}}$$

Attribute: Wind

$$S_{\text{strong}} = \{0+, 2-\} \Rightarrow E(S_{\text{strong}}) = \underline{\underline{0}}$$

$$S_{\text{weak}} = \{3+, 0-\} \Rightarrow E(S_{\text{weak}}) = \underline{\underline{0}}$$

$$\text{gain}(S_{\text{rain}}, \text{wind}) = \underline{\underline{0.97}}$$

Day	Temp	Humidity	Wind	Play
D4	Mild	high	weak	Yes
D5	Cool	normal	weak	Yes
D6	Cool	normal	Strong	No
D10	mild	normal	weak	Yes
D14	mild	high	Strong	No

$$\text{gain}(S_{\text{rain}}, \text{temp}) = \underline{\underline{0.019}}$$

$$\text{gain}(S_{\text{rain}}, \text{humid}) = \underline{\underline{0.019}}$$

$$\text{gain}(S_{\text{rain}}, \text{wind}) = \underline{\underline{0.97}}$$



wind \Rightarrow weak , strong

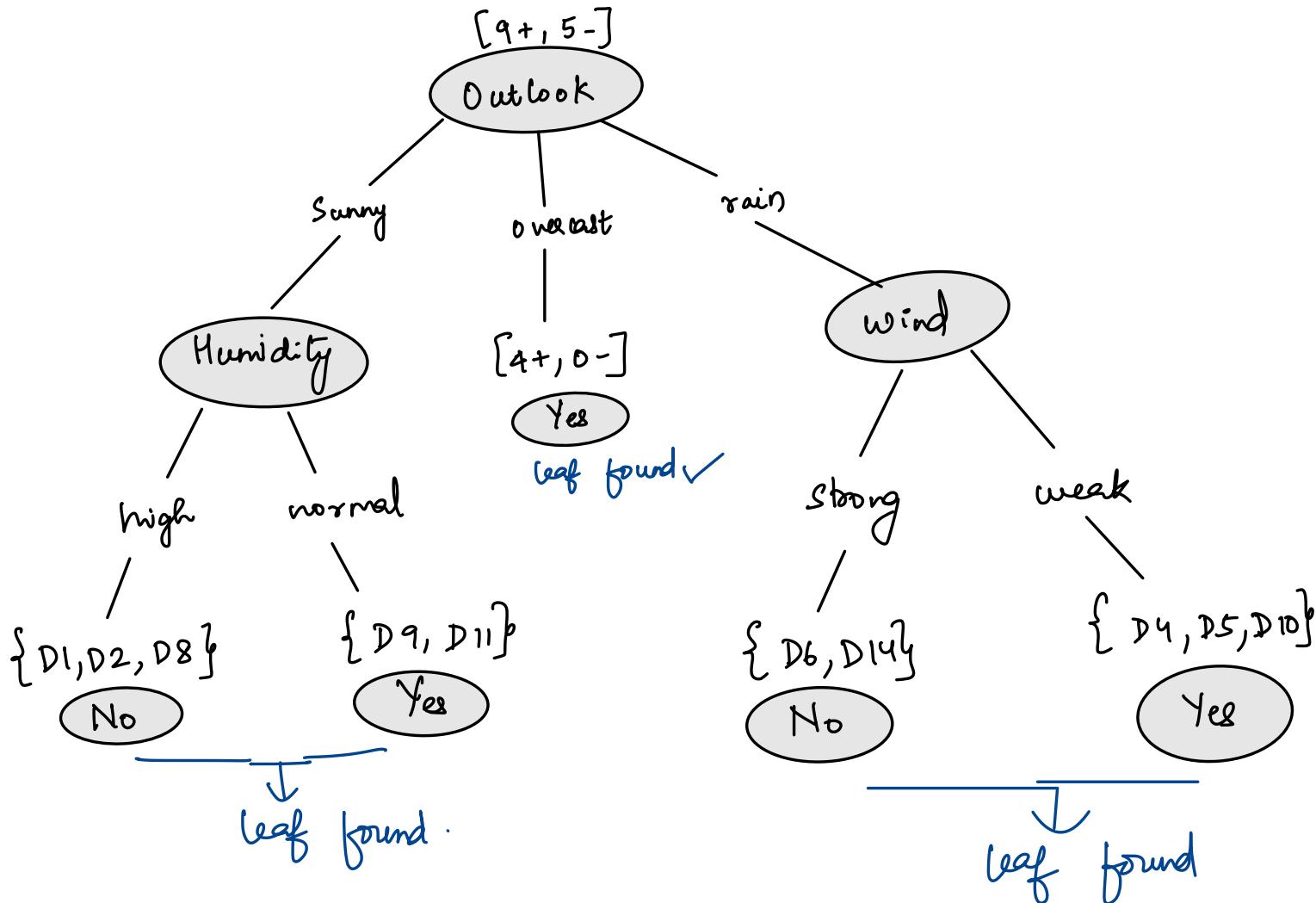


3 Yes

0 Yes

0 No

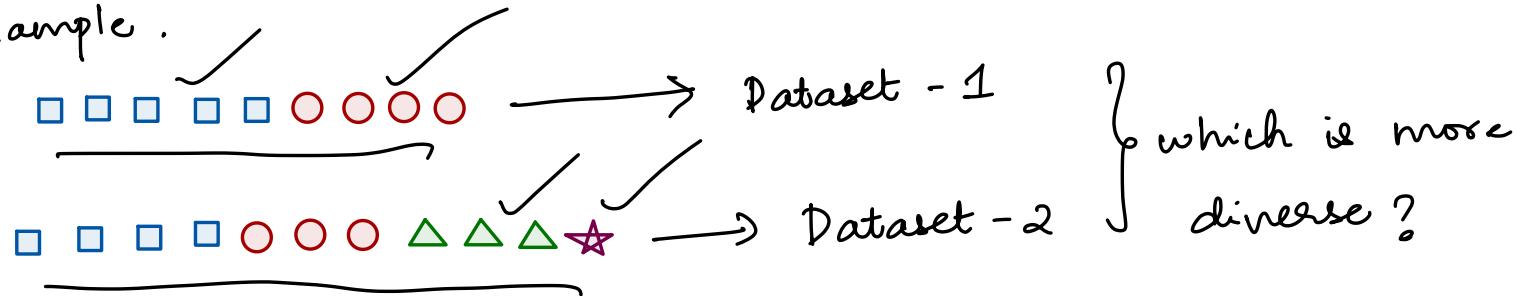
2 No



Gini Impurity Index

→ Measure of diversity in a dataset

Example.



Finding impurity:-

→ picking 2 random elements & check if they all same

<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Same
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Different
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Different
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Same
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Same
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Different
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Same
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Same
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Different
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Same

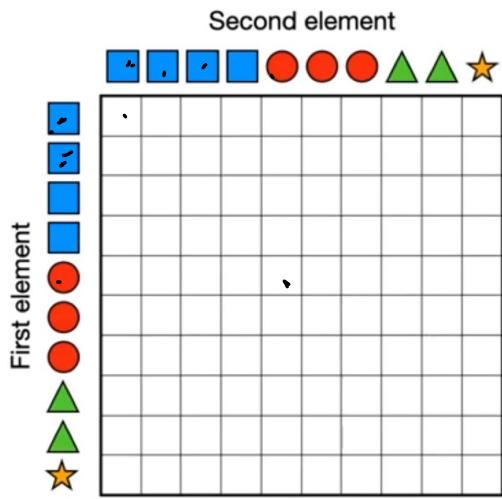
Different
= $\frac{4}{10}$
 4



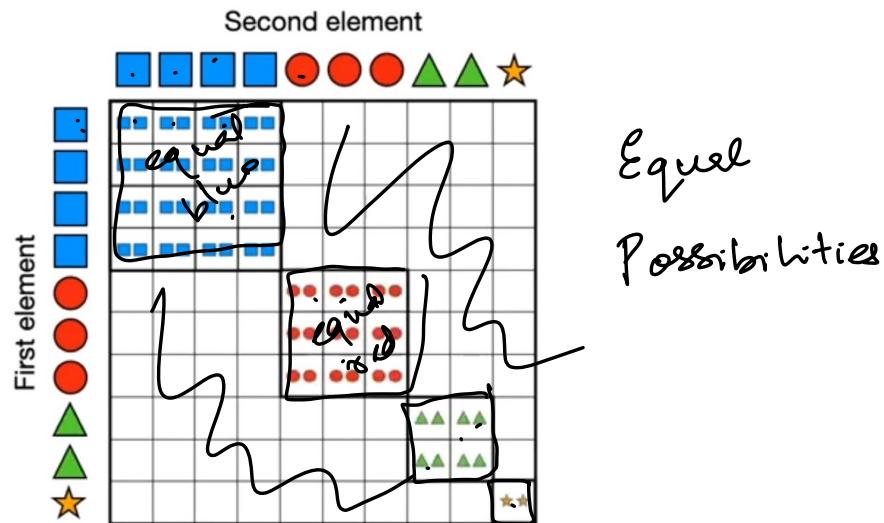
<input checked="" type="radio"/>	<input checked="" type="triangle"/>	Different
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Same
<input checked="" type="triangle"/>	<input checked="" type="radio"/>	Different
<input checked="" type="star"/>	<input checked="" type="radio"/>	Different
<input checked="" type="radio"/>	<input checked="" type="triangle"/>	Different
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Same
<input checked="" type="triangle"/>	<input checked="" type="radio"/>	Different
<input checked="" type="radio"/>	<input checked="" type="star"/>	Different
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Different

Different
= $\frac{7}{10}$
 7

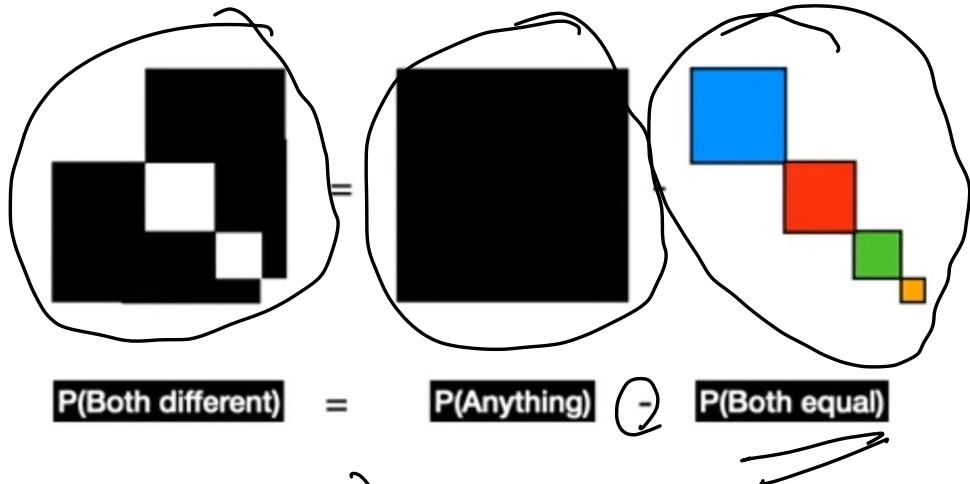
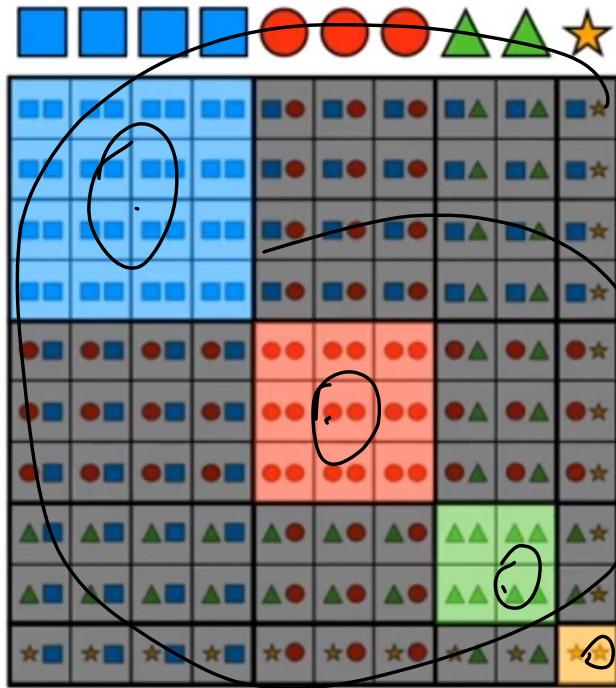
Bini önder = Probability of picking 2 distinct elements.



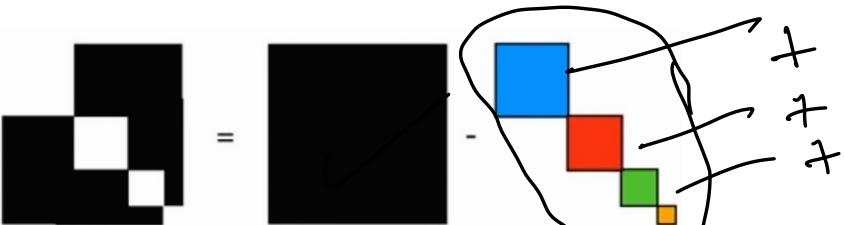
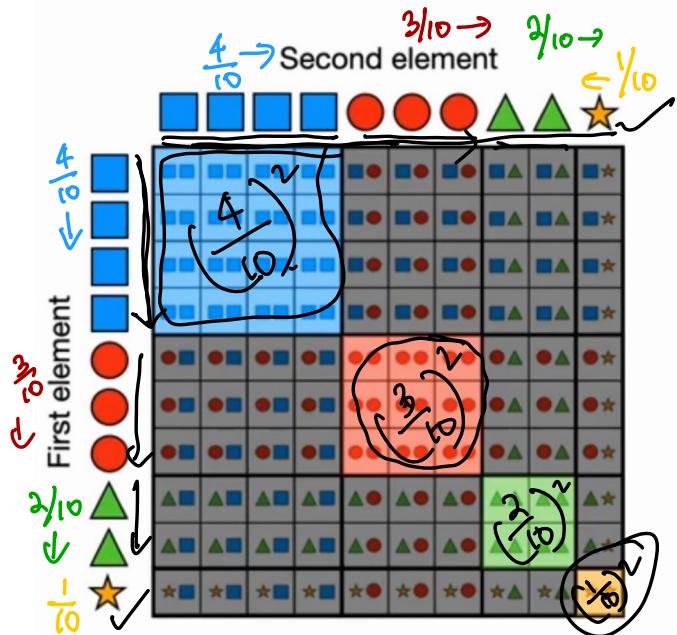
Considering entire data , all possibilities
 \checkmark
 $\frac{10}{2}$ elements
 $\frac{10}{2}$ elements



Second element



$$\frac{P(\text{anything})}{\Rightarrow 1} = \underline{\underline{1}}$$
$$\Rightarrow 1 = (P(\text{blue}) + P(\text{red}) + P(\text{green}) + P(\text{yellow}))$$

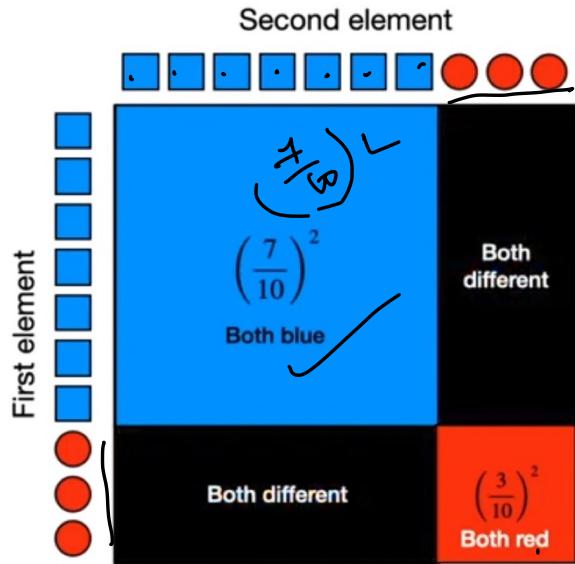


$$P(\text{Both different}) = P(\text{Anything}) - P(\text{Both equal})$$

$$= 1 - \left(\left(\frac{4}{10} \right)^2 + \left(\frac{3}{10} \right)^2 + \left(\frac{2}{10} \right)^2 + \left(\frac{1}{10} \right)^2 \right)$$

$$= 1 - (0.16 + 0.09 + 0.04 + 0.01)$$

$$P(\text{diff}) = 1 - 0.3 = 0.7$$



$$\begin{aligned}
 & P(\text{both different}) \\
 &= P(\text{Anything}) - (P(\text{both blue}) + \\
 &\quad P(\text{both red})) \\
 &= 1 - \left(\left(\frac{7}{10}\right)^2 + \left(\frac{3}{10}\right)^2 \right) \\
 &= 1 - (0.49 + 0.09) \\
 &= 1 - 0.58 = \underline{\underline{0.42}}
 \end{aligned}$$

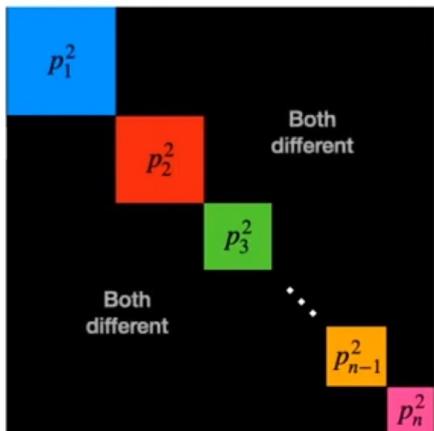
General formula

→ n classes

→ proportions → $p_1, p_2, p_3 \dots p_n$

→ Gini impurity index = $1 - (\text{sum of squares of proportions})$

$$1 - (p_1^2 + p_2^2 + \dots + p_n^2)$$



Some more examples

$$\overbrace{\square \square \square \square \circ \circ \circ} \rightarrow \underline{\underline{0.42}}$$

$$\overbrace{\square \square \square \square \circ \circ \circ \triangle \triangle \triangle \star} \rightarrow \underline{\underline{0.7}}$$

$$\begin{aligned} & \overbrace{\square \square \square \square \square \square \square \square \square} \rightarrow \underline{\underline{0}} \quad (1 \text{ (anything)}) \\ & = 1 - 1^2 = \underline{\underline{0}} \end{aligned}$$

$\left(\frac{1}{10}\right)^2$

$$\overbrace{\circ \circ \circ \circ \circ \circ \circ \circ \circ \circ} \rightarrow \text{all are different}$$

\Rightarrow every one will have $\left(\frac{1}{10}\right)^2$ also

$$\begin{aligned} & = 1 - \underbrace{0.1^2 + 0.1^2 + \dots + 0.1^2}_{10 \text{ times}} \\ & = \underline{\underline{0.9}} \end{aligned}$$

Example Dataset

Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay In
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

Step 1 → Compute Gini Index for entire dataset
 possible values of target
 ⇒ Cinema, Tennis, Stay in, Shopping (4) (total 10)

$$\rightarrow \text{Cinema} \Rightarrow \frac{6}{10}$$

$$\rightarrow \text{Tennis} \Rightarrow \frac{2}{10}$$

$$\rightarrow \text{Stay in} \Rightarrow \frac{1}{10}$$

$$\rightarrow \text{Shopping} \Rightarrow \frac{1}{10}$$

$$Gini(S) = 1 - \left(\left(\frac{6}{10}\right)^2 + \left(\frac{2}{10}\right)^2 + \left(\frac{1}{10}\right)^2 + \left(\frac{1}{10}\right)^2 \right) = \underline{\underline{0.58}}$$

Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay In
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

→ to get root node ⇒ compute
Gini of all attributes

① Weather → Sunny, windy, rainy

Sunny (3) ⇒ 1 Cinema, 2 Tennis

$$\text{Gini}(\text{Sunny}) = 1 - \left[\left(\frac{1}{3} \right)^2 + \left(\frac{2}{3} \right)^2 \right]$$

$$= 1 - [0.111 + 0.444] \approx \underline{\underline{0.44}}$$

Gini (S windy) :-

total - 4 ⇒ 3 cinema, 1 shopping

$$\Rightarrow 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 1 - [0.5625 + 0.0625] = \underline{\underline{0.375}}$$

Gini (S rainy) :-

total = 3 ⇒ 2 cinema, 1 stay in

$$\Rightarrow 1 - \left[\left(\frac{2}{3} \right)^2 + \left(\frac{1}{3} \right)^2 \right] = \underline{\underline{0.44}}$$

Final Gini: weighted avg

$$0.44 \times \left(\frac{3}{10} \right) + 0.375 \times \left(\frac{4}{10} \right)$$

$$+ 0.44 \times \left(\frac{3}{10} \right) \Rightarrow \underline{\underline{0.414}}$$

Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay In
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

② Parents \Rightarrow Yes, no

$$\text{gini } (S_{\text{Yes}}) \Rightarrow \underline{\text{total 5}}$$

\Rightarrow 5 cinema targets

$$\Rightarrow 1 - \left[\frac{5}{5} \right]^2 = 1 - 1^2 = \underline{\underline{0}}$$

$$\text{gini } (S_{\text{no}}) \Rightarrow \underline{\text{5 total}}$$

\Rightarrow 2 tennis, 1 stay in,
1 cinema, 1 shopping

$$\Rightarrow 1 - \left[\left(\frac{2}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 \right]$$

$$= 1 - [0.16 + 0.04 + 0.04 + 0.04] = \underline{\underline{0.72}}$$

Weighted average:-

$$0 * \frac{5}{10} + 0.72 * \frac{5}{10} \Rightarrow \underline{\underline{0.36}}$$

Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay In
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

Repeat for money attribute.
 weighted average of money = 0.486

Now compare all gini indexes:-

$$\text{gini}(\text{weather}) = 0.414$$

$$\text{gini}(\text{parents}) = \textcircled{0.36}$$

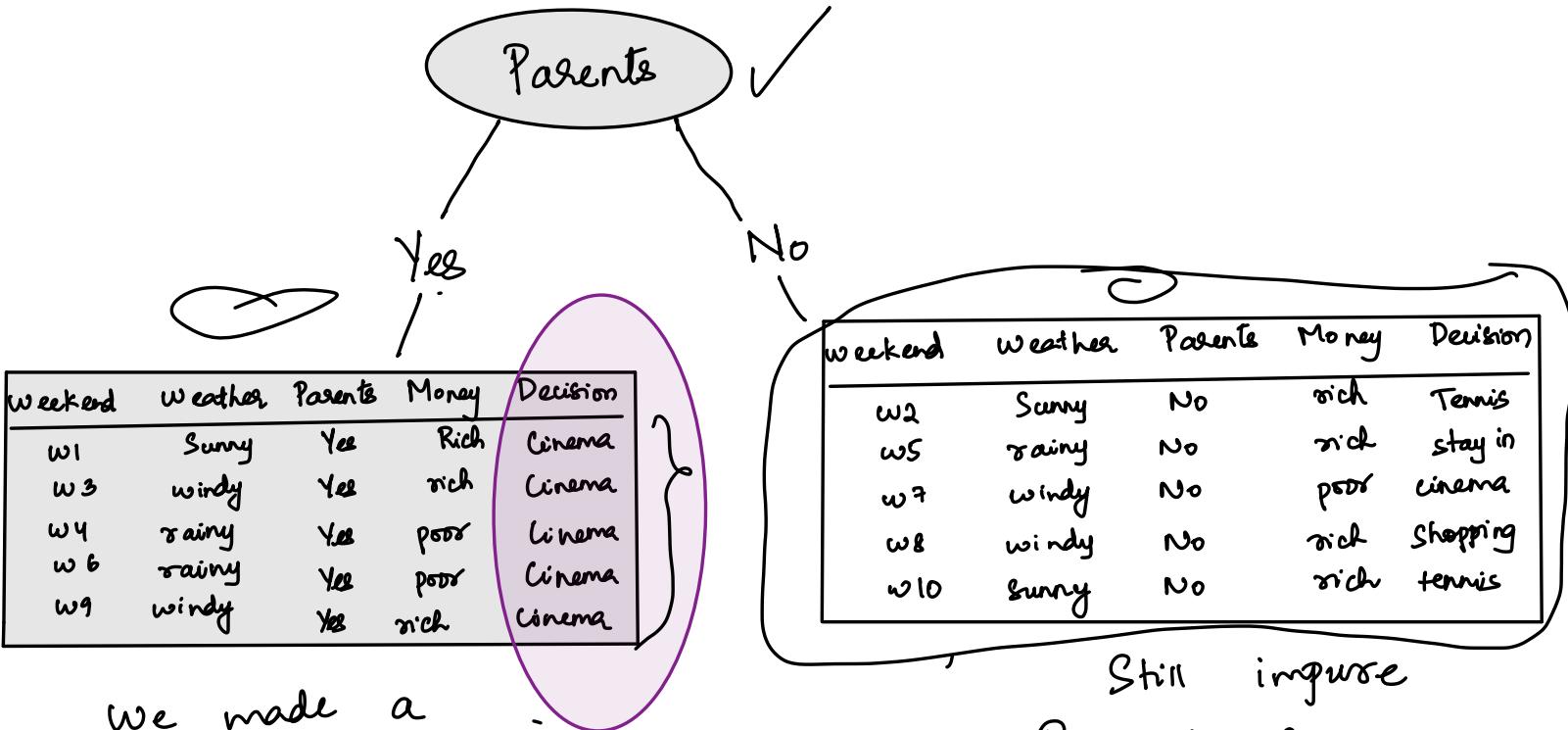
$$\text{gini}(\text{money}) = 0.486$$

Select the attribute with

least gini impurity

→ parents selected as root node.

→ from root, split the dataset with possible values
 ⇒ (here, yes or no)



We made a
decision here

Still impure

Repeat the
same process for
this sample

weekend	weather	Parents	Money	Decision
w2	Sunny	No	rich	Tennis
w5	rainy	No	rich	stay in
w7	windy	No	poor	cinema
w8	windy	No	rich	Shopping
w10	sunny	No	rich	tennis

$yini(S_{\text{weather}}) \Rightarrow \underline{\text{sunny}}, \underline{\text{rainy}}, \underline{\text{windy}}$

$yini(S_{\text{sunny}}) \Rightarrow \text{total } 2 ; \text{ both tennis} \Rightarrow \underline{\underline{0}}$

$yini(S_{\text{rainy}}) = 1 \text{ value} \Rightarrow \text{stay in} \Rightarrow \underline{\underline{0}}$

$yini(S_{\text{windy}}) = 2 \text{ values} \Rightarrow \text{cinema, shopping}$

$$= 1 - \left[\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right] = \underline{\underline{0.5}}$$

Weighted average = $0 \times \frac{2}{5} + 0 \times \frac{1}{5} + 0.5 \times \frac{2}{5} = \underline{\underline{0.2}}$

weekend	weather	Parents	Money	Decision
w2	Sunny	No	rich	Tennis
w5	rainy	No	rich	stay in
w7	windy	No	poor	cinema
w8	windy	No	rich	Shopping
w10	sunny	No	rich	tennis

$\text{gini}(\leq \text{money}) \Rightarrow \text{rich, poor}$

$\text{gini}(\leq \text{rich}) \Rightarrow \text{total 4 values}$

2 tennis, 1 stay in, 1 shopping

$$\Rightarrow 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{1}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = \underline{\underline{0.625}}$$

$\text{gini}(\leq \text{poor}) \Rightarrow 1 \text{ value} \Rightarrow \text{cinema} \Rightarrow \underline{0}$

Weighted average = $0.625 \times \frac{4}{5} + 0 \times \frac{1}{5} = \underline{\underline{0.5}}$

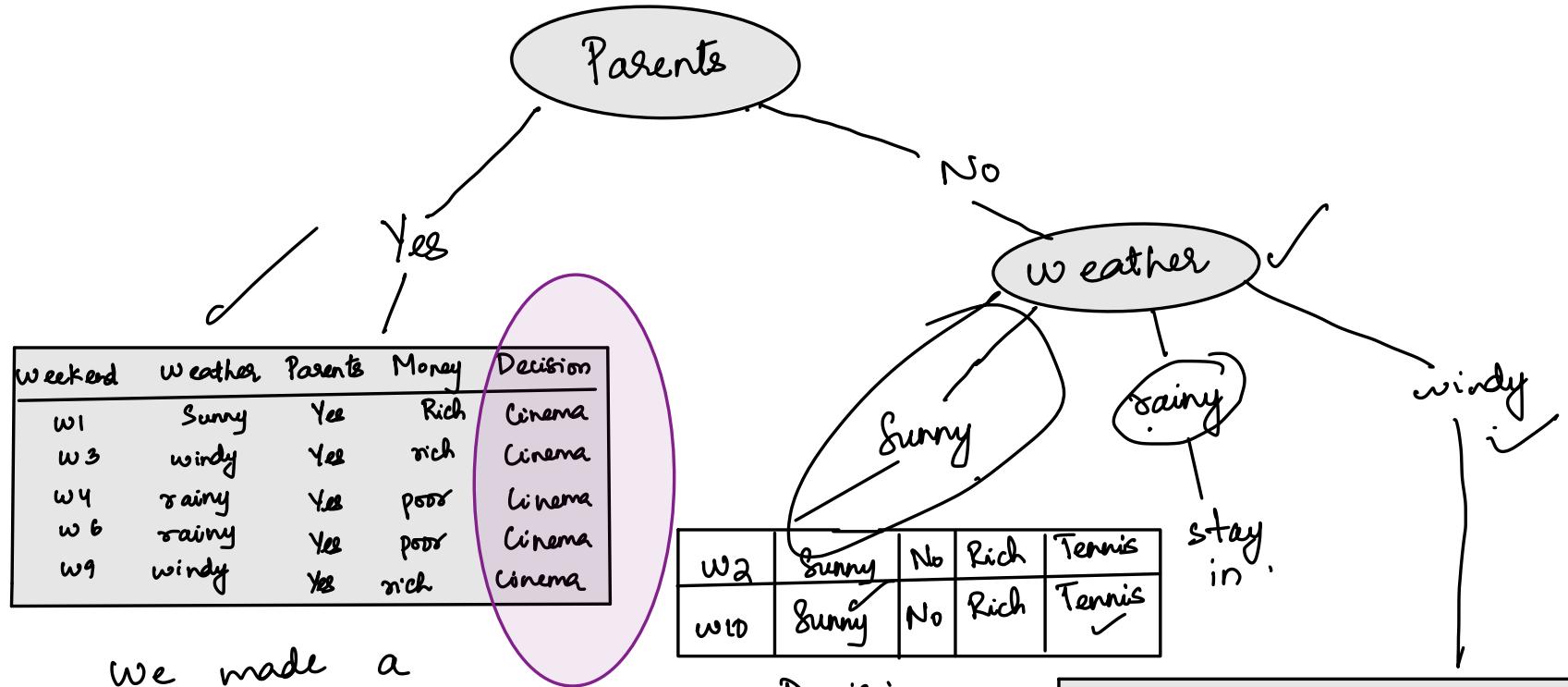
weekend	Weather	Parents	Money	Decision
w2	Sunny	No	rich	Tennis
w5	rainy	No	rich	stay in
w7	windy	No	poor	Cinema
w8	windy	No	rich	Shopping
w10	Sunny	No	rich	Tennis

Comparing the importances :-

$$\text{Weather} = 0.2$$

$$\text{Money} = 0.5$$

\Rightarrow Selected



We made a
decision here

Decision
made

w7 Windy No Poor Cinema
w8 Windy No Rich Shopping

Ned further
decision

weekend	weather	Parents	Money	Decision
w7	windy	no	poor	Cinema ✓
w8	windy	no	rich	Shopping ✓

Only 2 values for money

rich & poor

can directly

build a tree

\Rightarrow poor \Rightarrow cinema

\Rightarrow rich \Rightarrow shopping

Final tree

