

ID3 Algorithm Manually Worked Through

Play Outdoor Sport

Num	Outlook	Temperature	Humidity	Wind	Play?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cold	Normal	Weak	Yes
6	Rain	Cold	Normal	Strong	No
7	Overcast	Cold	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cold	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Jared O'Toole

2019-11-18

Oracle Machine Learning
Section 4-3

① Calculate Entropy on the results column "Play"

- there are 9 "yes" and 5 "no" (2 outcomes)
- $\text{Entropy}(S) = -P_{\text{yes}} \log_2(P_{\text{yes}}) + -P_{\text{no}} \log_2(P_{\text{no}})$
- $= -\frac{9}{14} \log_2\left(\frac{9}{14}\right) + -\frac{5}{14} \log_2\left(\frac{5}{14}\right)$
- ≈ 0.94

② Find the next attribute with the most gain

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{values}(A)} \left[\frac{|S_v|}{|S|} \text{Entropy}(S_v) \right]$$

③ \Rightarrow Test the "Outlook" Column

- There are 3 options: "Sunny", "Overcast", and "Rain" Total: 14

Outlook	# with "Yes" result	# with "No" result	Total #
Sunny	11	2	13
Overcast	4	0	4
Rain	3	11	14

$$\bullet \text{Entropy}(\text{Sunny}) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) + -\frac{3}{5} \log_2\left(\frac{3}{5}\right) \approx 0.971$$

$$\bullet \text{Entropy}(\text{Overcast}) = -\frac{4}{4} \log_2\left(\frac{4}{4}\right) + -\frac{0}{4} \log_2\left(\frac{0}{4}\right) = 0$$

$$\bullet \text{Entropy}(\text{Rain}) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) + -\frac{2}{5} \log_2\left(\frac{2}{5}\right) \approx 0.971$$

$$\bullet \text{Gain}(S, \text{Outlook}) = 0.94 - \left(V_{E_{\text{sunny}}} + V_{E_{\text{overcast}}} + V_{E_{\text{rain}}} \right)$$

with $V_E = \frac{|S_v|}{|S|} \cdot \text{Entropy}(S_v)$

$$\bullet V_{E_{\text{sunny}}} = (5/14)(0.971) = 0.357$$

$$\bullet V_{E_{\text{overcast}}} = (4/14)(0) = 0.000$$

$$\bullet V_{E_{\text{rain}}} = (5/14)(0.971) = 0.357$$

$$\dots = 0.94 - 0.357 - 0.357 = 0.246$$

(b) \Rightarrow Test the "Temperature" Column

- There are 3 options: "Cold", "mild", and "hot" Total: 14

Temperature	# with "Yes" result	# with "No" result	Total #
Cold	3	1	4
Mild	4	2	6
Hot	2	2	4

$$\bullet \text{Entropy}(\text{cold}) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) + -\frac{1}{4} \log_2\left(\frac{1}{4}\right) = 0.8112781245\dots$$

$$\bullet \text{Entropy}(\text{mild}) = -\frac{4}{6} \log_2\left(\frac{4}{6}\right) + -\frac{2}{6} \log_2\left(\frac{2}{6}\right) = 0.9182958341\dots$$

$$\bullet \text{Entropy}(\text{hot}) = -\frac{2}{4} \log_2\left(\frac{2}{4}\right) + -\frac{2}{4} \log_2\left(\frac{2}{4}\right) = 1.0$$

$$\bullet V_{E_{\text{cold}}} = (4/14)(0.8112781245\dots) = 0.2317937499\dots$$

$$\bullet V_{E_{\text{mild}}} = (6/14)(0.9182958341\dots) = 0.3935553575\dots$$

$$\bullet V_{E_{\text{hot}}} = (4/14)(1) = 0.2857142857$$

$$\bullet \text{Gain}(S, \text{Temperature}) = 0.94 - (V_{E_{\text{cold}}} + V_{E_{\text{mild}}} + V_{E_{\text{hot}}})$$

$$= 0.0289366069\dots$$

(c) \Rightarrow Test the "Humidity" Column

- There are 2 options: "high" and "normal" Total: 14

Humidity	# with "Yes" result	# with "No" result	Total #
High	III	3	7
Normal	NNI	6	7

- Entropy_S(high) = $-\frac{3}{7} \log_2\left(\frac{3}{7}\right) - \frac{4}{7} \log_2\left(\frac{4}{7}\right) = 0.985228136\dots$
- Entropy_S(normal) = $-\frac{6}{7} \log_2\left(\frac{6}{7}\right) - \frac{1}{7} \log_2\left(\frac{1}{7}\right) = 0.5916727786\dots$
- $V_{E_{\text{high}}} = (7/14)(0.985228136\dots) = 0.462614068\dots$
- $V_{E_{\text{normal}}} = (7/14)(0.5916727786\dots) = 0.2958363893\dots$
- Gain(S, Humidity) = $0.94 - (V_{E_{\text{high}}} + V_{E_{\text{normal}}})$
 $= 0.1515495427\dots$

(d) ▷ Test the "Wind" column

- There are 2 options: "Strong" and "Weak" Total: 14

Wind	# with "Yes" result	# with "No" result	Total #
Strong	III	5	8
Weak	III	3	6

- Entropy_S(Strong) = $-\frac{6}{8} \log_2\left(\frac{6}{8}\right) - \frac{2}{8} \log_2\left(\frac{2}{8}\right) = 0.8112781245\dots$
- Entropy_S(Weak) = $-\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right) = 1.0$
- $V_{E_{\text{strong}}} = (8/14)(0.8112781245\dots) = 0.4635874997\dots$
- $V_{E_{\text{weak}}} = (6/14)(1) = 0.4285714286\dots$
- Gain(S, Wind) = $0.94 - (V_{E_{\text{strong}}} + V_{E_{\text{weak}}})$
 $= 0.0478410717\dots$

(e) ▷ choose the column with the highest gain to be the root node

• Gain(S, Outlook) = 0.246 → Outlook

• Gain(S, Temperature) = 0.0289366069...

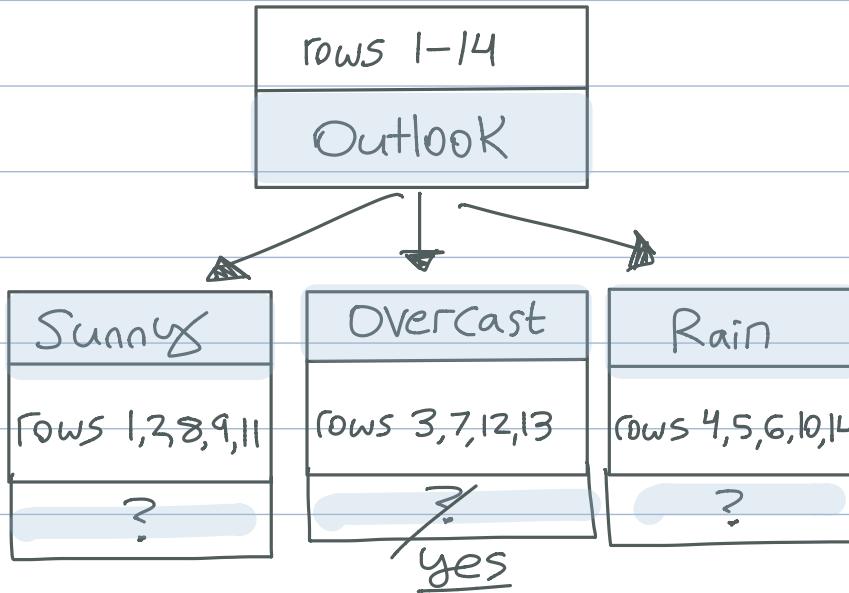
• Gain(S, Humidity) = 0.1515495427...

$$\cdot \text{Gain}(s, \text{Wind}) = 0.0478410717\dots$$

⇒ g h i

(f)

⇒ Append the tree



• "Sunny" node has $\frac{2}{5}$ "yes's

$$E(\text{Sunny}) = 0.9709505945$$

$$\left(= -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right)\right)$$

• "Overcast" node has $\frac{3}{4}$ "yes's

It can be categorized as "Yes"

• "Rain" node has $\frac{3}{5}$ "Yes's

$$E(\text{Rain}) = 0.9709505945$$

$$\left(= -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right)\right)$$

③ Find the next attribute with the most gain to split the "Sunny" node

⇒ a b c Test the "Temperature", "Humidity", and "Wind" columns Total: 5 rows

Temp	# Yes	Totals
hot	0	2
mild	1	2
cold	1	1

Humid.	# Yes	Totals
High	0	3
Normal	2	2

Wind	# Yes	Totals
weak	1	3
Strong	1	2

$$E(\text{hot}) = 0$$

$$\left(= -\frac{0}{2} \log_2 \left(\frac{0}{2}\right) - \frac{2}{2} \log_2 \left(\frac{2}{2}\right)\right)$$

$$E(\text{mild}) = 1$$

$$\left(= -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right)$$

$$E(\text{cold}) = 0$$

$$\left(= -\frac{1}{1} \log_2 \left(\frac{1}{1}\right) - \frac{0}{1} \log_2 \left(\frac{0}{1}\right)\right)$$

$$VE_{\text{hot}} = \frac{2}{5} \cdot 0 = 0$$

$$VE_{\text{mild}} = \frac{2}{5} \cdot 1 = 0.40$$

$$VE_{\text{cold}} = \frac{1}{5} \cdot 0 = 0$$

$$\text{Gain}(\text{Temp.}) = 0.57095\dots$$

$$E(\text{high}) = 0$$

$$E(\text{normal}) = 0$$

$$VE_{\text{high}} = 0$$

$$VE_{\text{normal}} = 0$$

$$\text{Gain}(\text{Humid.}) = 0.9709\dots$$

$$\left(= 0.97 - (0+0)\right)$$

$$E(\text{weak}) = 0.9182958341\dots$$

$$\left(= -\frac{1}{3} \log_2 \left(\frac{1}{3}\right) - \frac{2}{3} \log_2 \left(\frac{2}{3}\right)\right)$$

$$E(\text{strong}) = 1.0$$

$$VE_{\text{weak}} = \frac{3}{5} \cdot E(\text{weak}) = 0.5509775\dots$$

$$VE_{\text{strong}} = \frac{2}{5} \cdot E(\text{strong}) = 0.40$$

$$\text{Gain}(\text{Wind}) = 0.019973094$$

$$\left(= 0.97 - (0.551 + 0.4)\right)$$

$$= 0.971 - (V_{E\text{hot}} + V_{E\text{mild}} + V_{E\text{cold}})$$

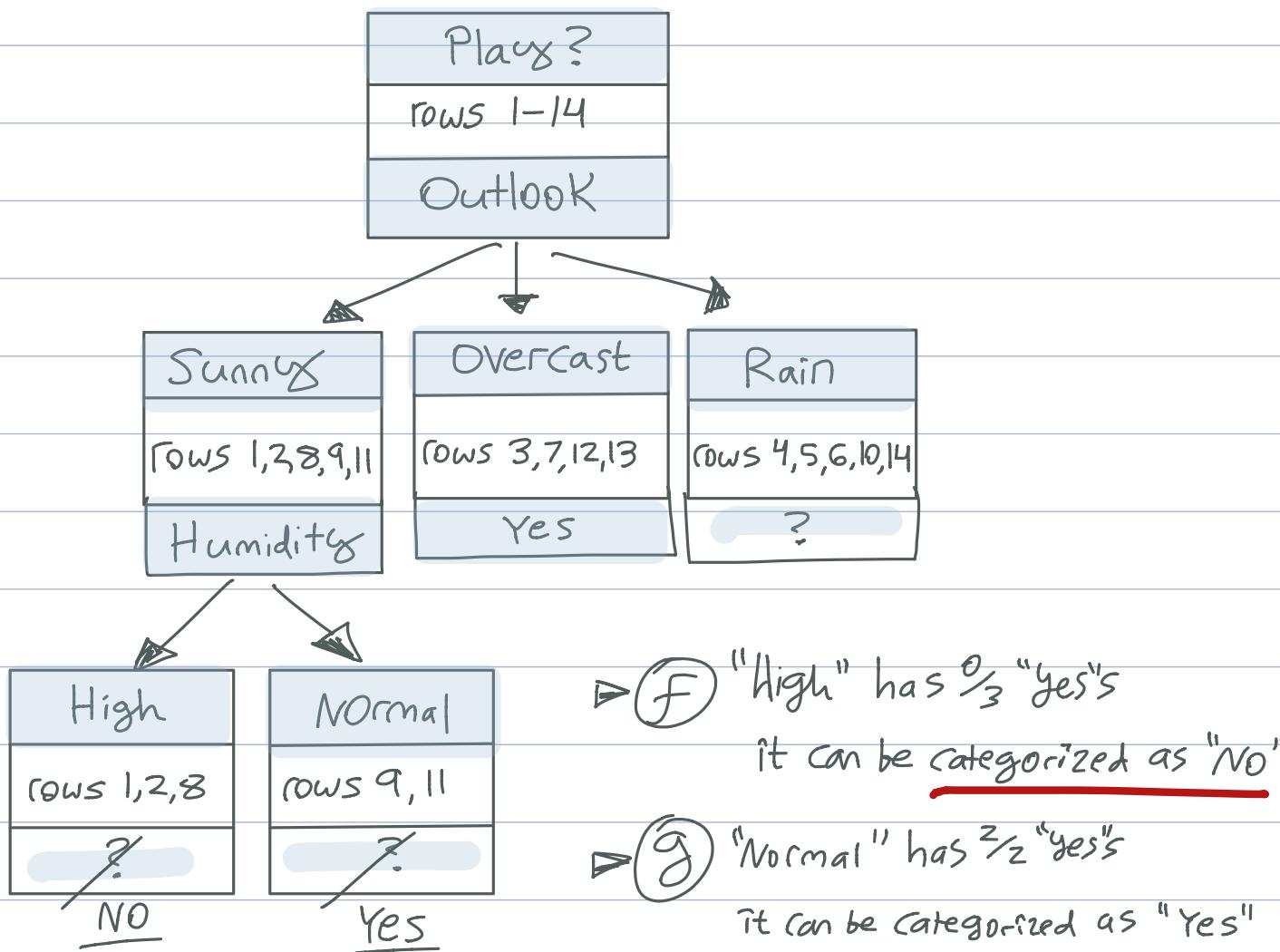
⇒ (d) Choose the column with the highest Gain
to decide the child nodes of "Sunny"

- Gain(Sunny, Temp) = 0.57095...

- Gain(Sunny, Humid) = 0.9709... → **Humidity**

- Gain(Sunny, Wind) = 0.019973094

⇒ (e) Append the tree



- (4) Find the next attribute with the most gain to split the "Rain" node
- ⇒ (a)(b)(c) Test the "Temperature", "Humidity", and "Wind" columns

Total: 5 rows

Temp	# Yes	Totals
hot	—	0
mild	2	3
cold	1	2

Humid.	# Yes	Totals
High	1	2
Normal	2	3

Wind	# Yes	Totals
weak	3	3
strong	0	2

$$E(\text{high}) = \underline{1}$$

$$E(\text{normal}) = \underline{0.918\dots}$$

$$E(\text{mild}) = \underline{0.9182958341\dots}$$

$$E(\text{cold}) = \underline{1}$$

$$DE_{\text{high}} = \frac{2}{5} \cdot 1 = \underline{0.4}$$

$$DE_{\text{mild}} = \frac{3}{5} E(\text{mild}) = \underline{0.55097\dots}$$

$$DE_{\text{norm}} = \frac{3}{5} \cdot E(\text{norm}) = \underline{0.5509\dots}$$

$$DE_{\text{cold}} = \frac{2}{5}(1) = \underline{0.40}$$

$$\text{Gain}_{\text{Humid}}^{(\text{Rain},)} = \underline{0.019973\dots}$$

$$\text{Gain}_{\text{Temp}}^{(\text{Rain},)} = \underline{0.019973094\dots}$$

$$E(\text{weak}) = \underline{0}$$

$$E(\text{strong}) = \underline{0}$$

$$DE_{\text{weak}} = \frac{3}{5} \cdot 0 = \underline{0}$$

$$DE_{\text{strong}} = \frac{2}{5} \cdot 0 = \underline{0}$$

$$\text{Gain}_{\text{Wind}}^{(\text{Rain},)} = \underline{0.9709\dots}$$

⇒ (d) Choose the column with the highest gain for "Rain"

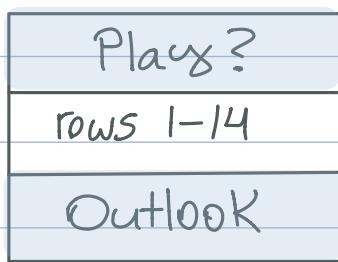
• $\text{Gain}(\text{Wind}, \text{Temp}) = \text{Gain}(\text{Wind}, \text{Humid}) = 0.019973094$

• $\text{Gain}(\text{Wind}, \text{Wind}) = 0.9709\dots \rightarrow \boxed{\text{Wind}}$

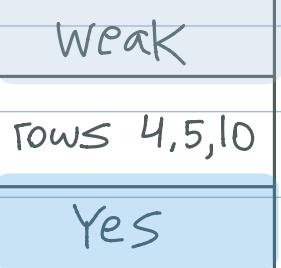
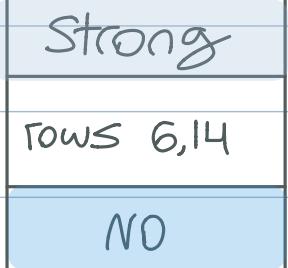
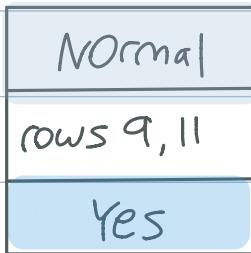
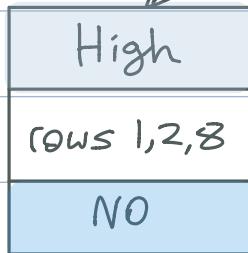
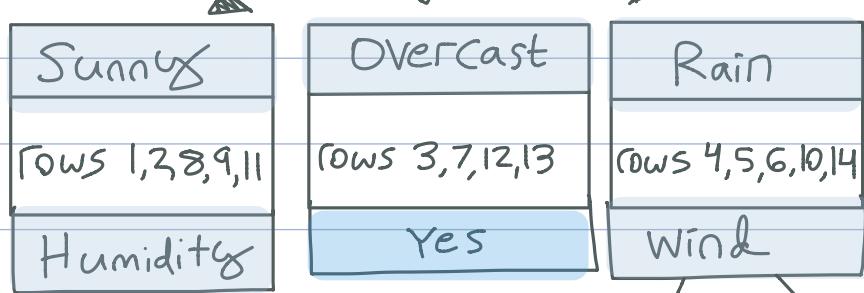
⇒ (e) Append the tree

⇒ (f) "strong" has 0/2 "yes's
it can be classified "No"

The final tree:



⇒ (g) "weak" has 3/3 "yes's
it can be classified "Yes"



Pseudo-code for the ID3 algorithm:

```
# RECURSIVE FUNCTION

def branch(node, table)

    px = table.successes / table.total

    if px in (0, 1) # BASE CASE

        node.result = px

        return

    if not table.attrs # BASE CASE

        node.result = bool math.round(px)

        return

    # RECURSIVE CASE

    entropy_s = entropy(table)

    gains = dict()

    for attr in table.attrs

        gains[attr] = gain(entropy_s, attr)

    best, max = null, 0

    for attr, g in gains.items()

        if g > max

            best, max = attr, g

    sliced_table = table.without(best)

    for outcome in best

        sub_table = sliced_table.with_only(outcome)

        child = Node()

        node.add(child)

        branch(child, sub_table)
```

