# Decision Tree Construction!

**Step 1:**

| Play | |
|------|------|
| Yes | No |
| 9 | 5 |
| 9/14 | 5/14 |

Pr.

$$H(Play?) = -\frac{9}{14}\log_2\frac{9}{14} - \frac{5}{14}\log\frac{5}{14} = 0.94$$

Entropy is a measure of uncertainity

**Goal: Find the most simple decision tree**

Split the table (tree) according to the "attribute" with the highest "information gain".

$$G(X) \triangleq H(Play?) - H(X) = H(S) - H(X)$$

where $X$ represent an attribute.     Source, original table

**Step 2:** Calculate the Gain of each attribute.

**2a)** Let $X =$ Temperature. (which takes 3 values cool, mild, hot)

| T | Play | | | |
|---|---|---|---|---|
| | | Yes | No | Bernoulli: |
| | C | 3 | 1 | 3/4 1/4 |
| | M | 4 | 2 | 4/6 2/6 |
| | H | 2 | 2 | |

$$H(C) = -\frac{3}{4}\log\frac{3}{4} - \frac{1}{4}\log\frac{1}{4}$$

$$H(M) = -\frac{4}{6}\log\frac{4}{6} - \frac{2}{6}\log\frac{2}{6}$$

$$H(Hot) = -\frac{1}{2}\log\frac{1}{2} - \frac{1}{2}\log\frac{1}{2}$$

$$H(T) = \frac{4}{14}H(C) + \frac{6}{14}H(M) + \frac{4}{14}(H(Hot)) = 0.911$$

Weighted average of $H(C)$, $H(M)$ & $H(Hot)$

$$G(T) = 0.94 - 0.911 = 0.029 \qquad \text{(small gain)}$$

Because $H(T)$ is high (uncertainity)     $G(Outlook) = 0.246$

**2b)** Let $X =$ Outlook.     $H(Outlook) = \frac{5}{14}H(S) + 0 + \frac{5}{14}H(R) = 0.6....$

| | Play | | Pr. |
|---|---|---|---|
| Outlook | Yes | No | |
| S | 2 | 3 | 2/5, 3/5 |
| O | 4 | 0 | 1, 0 |
| R | 3 | 2 | 3/5, 2/5 |

$$H(S) = -\frac{2}{5}\log\frac{2}{5} + \left(\frac{-3}{5}\log\frac{3}{5}\right)$$
$$H(O) = 0 \qquad \text{overcast} = 0$$
$$H(R) = H(Sunny) = \left(-\frac{3}{5}\log\frac{3}{5} - \frac{2}{5}\log\frac{2}{5}\right)$$
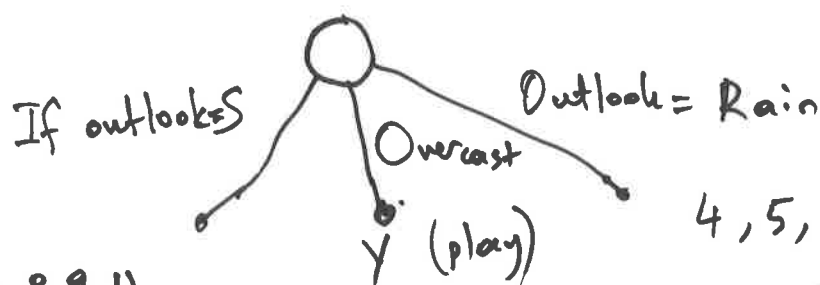
2 c) Let X= Humidity.

$$G(Humidity) = 0.151$$

2 d) Let X= Wind, $\quad G(Wind) = 0.048$

split the tree according to Outlook.

$$Outlook = \arg\max_{X} G(X)$$

Outlook is the attribute with the maximum information gain.

3)



If outlook=S

Overcast

Y (play)

Outlook = Rain

4, 5, 6, 10, 14 (rows)

rows: 1, 2, 8, 9, 11

| | H | W | T | |
|---|---|---|---|---|
| 1 | high | Weak | hot | N |
| 2 | H | H | S | N |
| 8 | | | | N |
| 9 | | | | Y |
| 11 | | | | Y |

| | H | W | T | Play |
|---|---|---|---|---|
| 4 | M | H | W | Y |
| 5 | C | N | W | Y |
| 6 | | | | N |
| 10 | | | | Y |
| 14 | | | | N |

$$H(ST1) = -\tfrac{3}{5}\log\tfrac{3}{5} - \tfrac{2}{5}\log\tfrac{2}{5}$$

$$H(ST2) = H(ST1)$$
.
.
.
.

Iteratively continue.

# Will I play tennis today?

| | O | T | H | W | Play? |
|---|---|---|---|---|---|
| ✓ 1 | S | H | H | W | - |
| ✓ 2 | S | H | H | S | - |
| 3 | O | H | H | W | + |
| 4 | R | M | H | W | + |
| ✓ 5 | R | C | N | W | + |
| 6 | R | C | N | S | - |
| ✓ 7 | O | C | N | S | + |
| ✓ 8 | S | M | H | W | - |
| ✓ 9 | S | C | N | W | + |
| 10 | R | M | N | W | + |
| ✓ 11 | S | M | N | S | + |
| 12 | O | M | H | S | + |
| 13 | O | H | N | W | + |
| 14 | R | M | H | S | - |

Entropy: Bernoulli distribution.

$P_+ = \frac{9}{14}$ , $P_- = \frac{5}{14}$

$H(Table) = -\frac{9}{14} \log \frac{9}{14} - \frac{5}{14} \log \frac{5}{14}$

"9 + & 5 -"

**Outlook:**
O(vercast),
S(unny),
R(ainy)

**Temperature:** H(ot),
M(edium),
C(ool)

**Humidity:** H(igh),
N(ormal),
L(ow)

**Wind:** S(trong),
W(eak)

$H(Table) = 0.94$
$H(Play?) = 0.94.$

# Example || Given the table

| A1 | A2 | Result | Pr |
|----|----|--------|-----|
| 0 | 0 | 0 | } ⟹ ½ |
| 0 | 1 | 0 | |
| 1 | 0 | 1 | } ⟹ ½ |
| 1 | 1 | 1 | |

$$H(R) = -\frac{1}{2}\log\frac{1}{2} - \frac{1}{2}\log\frac{1}{2}$$
$$H(R) = 1.$$

**Gain of A1:**

A1 "0" → "0"       Pr. 0, 1
A1 "1" → "1"       Pr. 1, 0

$$H(A1) = \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 0 = 0$$

(Entropy)

| A1 | Yes | No | Pr. |
|----|-----|-----|------|
| 0 | 0 | 2 | 0, 1 |
| 1 | 2 | 0 | 1, 0 |

$$H("0") = -0\log 0 - 1\log 1 = 0$$
$$H("1") = -1\log 1 - 0\log 0 = 0$$

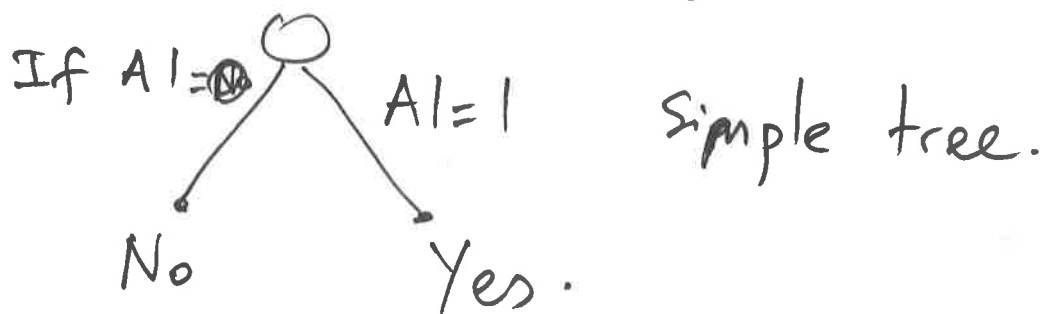$$G(A1) = 1 = 1 - 0. \quad \text{(very high)} \qquad (1*)$$

**Gain of A2:**

$$H(A_2) = \frac{1}{2}\underbrace{\left(-\frac{1}{2}\log\frac{1}{2} - \frac{1}{2}\log\frac{1}{2}\right)}_{H("0")''} + \frac{1}{2}\underbrace{\left(-\frac{1}{2}\log\frac{1}{2} - \frac{1}{2}\log\frac{1}{2}\right)}_{H("1")''}$$

| A2 | Yes | No | Pr. |
|----|-----|-----|------|
| 0 | 1 | 1 | ½, ½ |
| 1 | 1 | 1 | ½, ½ |

$$H(A_2) = 1$$

$$G(A2) = 0 = 1 - 1 = 0 \quad (\text{low}) \;②*$$

From (1*) & (2*) split according to A1.
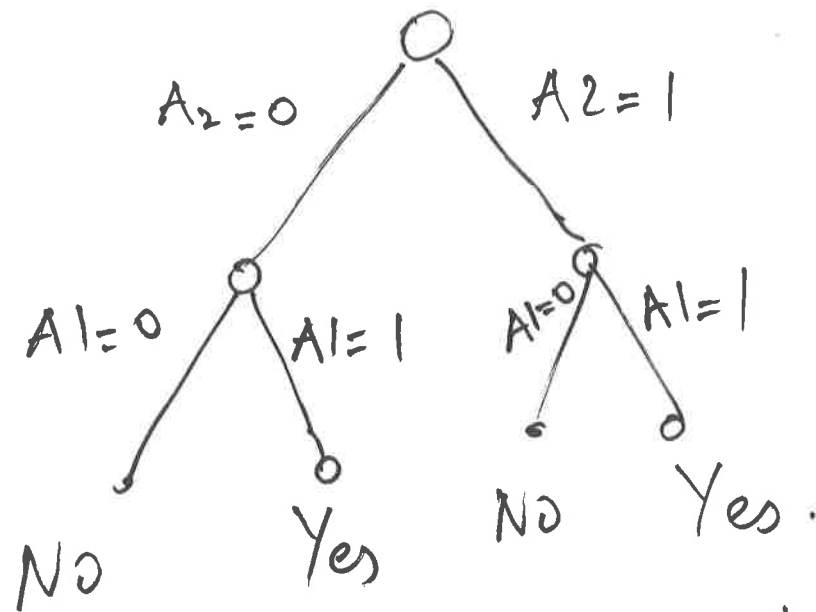


If A1=0          A1=1        Simple tree.

No          Yes.

At the leaves of the tree there is no
uncertainity. That is why we try to reduce
entropy.

# Inefficient tree:

Split
w.r.t A2

$A_2 = 0$      $A2 = 1$

$Tree\ 2$

$A1 = 0$   $A1 = 1$   $A1 = 0$   $A1 = 1$

No      Yes      No      Yes.

Computational complexity of the
tree is determined by the ✳ of comparisons
(questions) that you have to make (evaluate).

∴ Tree 2 is inefficient.