# Practical Machine Learning
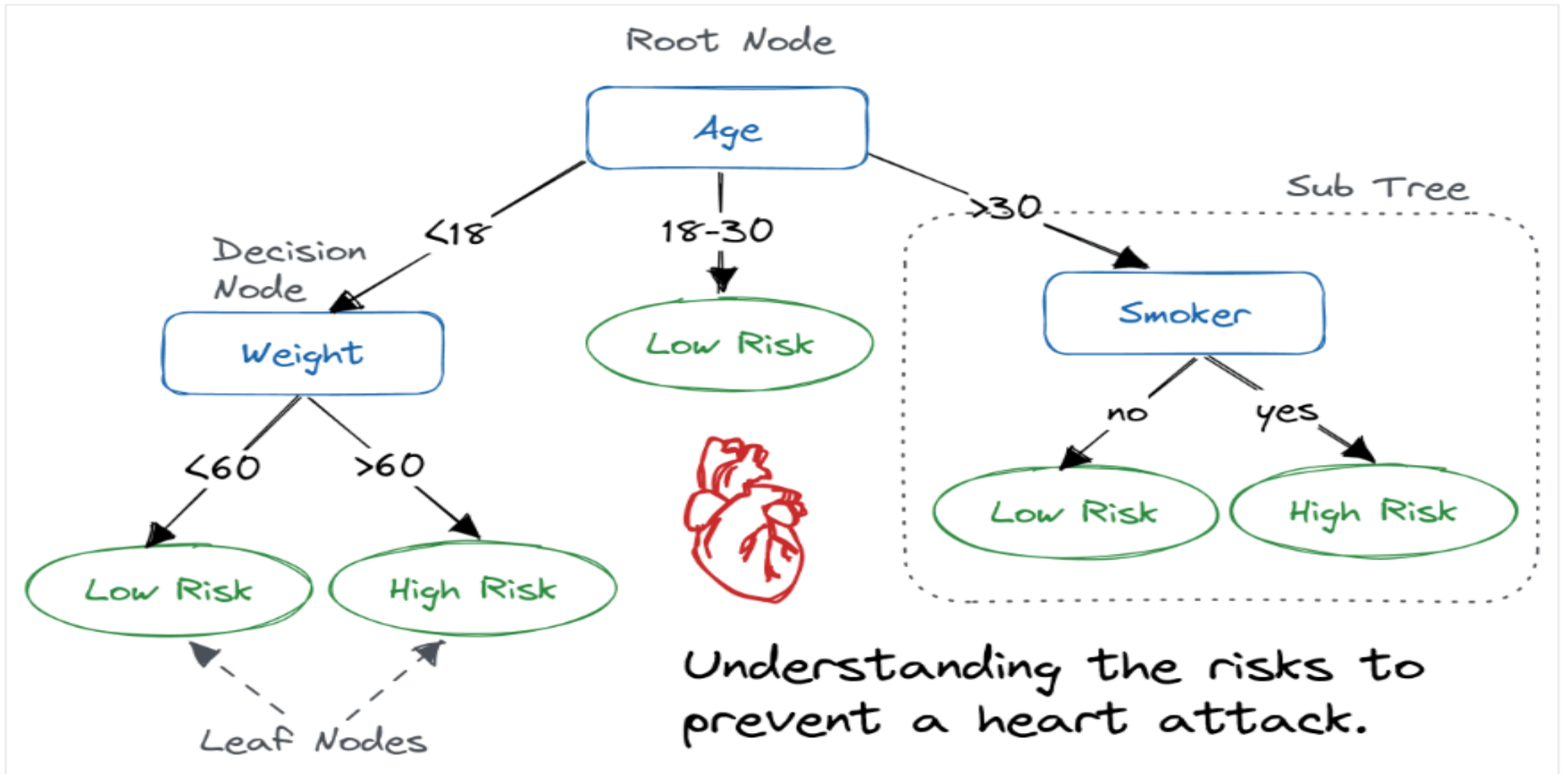
# Day 7: SEP23 DBDA

Kiran Waghmare

# Agenda

- Decision Tree
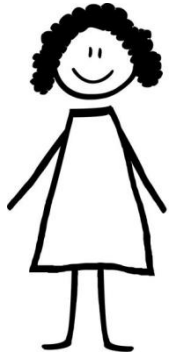
# Example :



Understanding the risks to prevent a heart attack.

# Definition

- A tree-like model that illustrates series of events leading to certain decisions
- Each node represents a test on an attribute and each branch is an outcome of that test

## Who to loan?

- Not a student
- 45 years old
- Medium income
- Fair credit record
- Student
  ➤ Yes
- 27 years old
- Low income
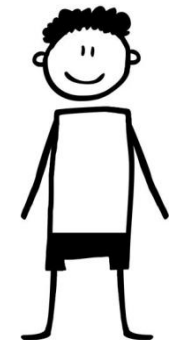- Excellent credit

# Definition

- A tree-like model that illustrates series of events leading to certain decisions
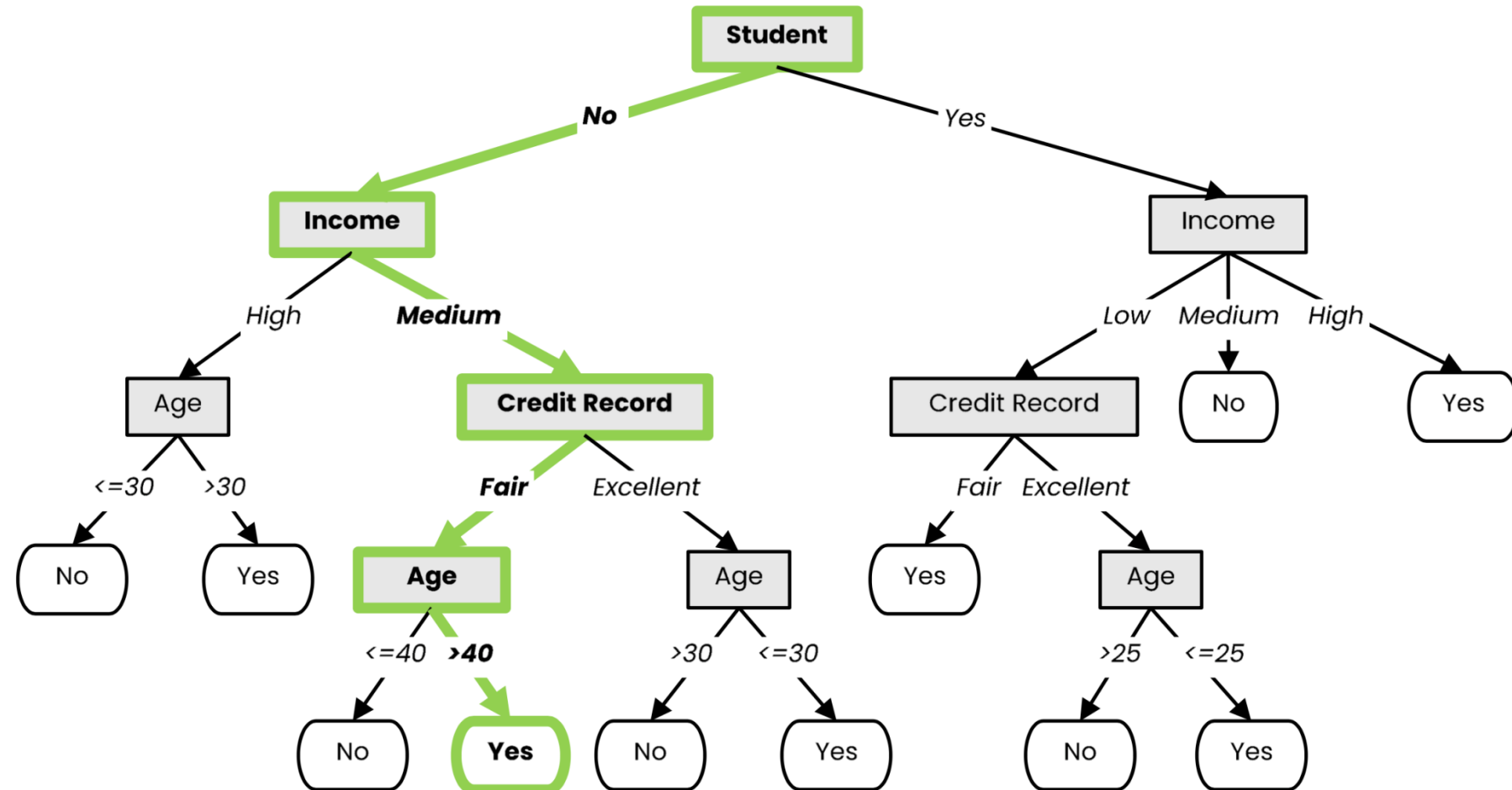- Each node represents a test on an attribute and each branch is an outcome of that test

Who to loan?

- Not a student
- 45 years old
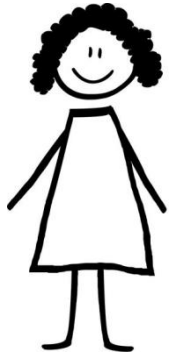- Medium income
- Fair credit record
- Student
  ➢ Yes
- 27 years old
- Low income
- Excellent credit

# Decision Tree

Root node

Branch

Decision node

Decision node

Leaf node

Decision node

Leaf node

Leaf node

Leaf node

Leaf node

Maximum Depth

# Process (1): Model Construction



| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Mike | Assistant Prof | 3 | no |
| Mary | Assistant Prof | 7 | yes |
| Bill | Professor | 2 | yes |
| Jim | Associate Prof | 7 | yes |
| Dave | Assistant Prof | 6 | no |
| Anne | Associate Prof | 3 | no |

Training Data

Classification Algorithms

Classifier (Model)

IF rank = 'professor'
OR years > 6
THEN tenured = 'yes'

# Process (2): Using the Model in Prediction

Classifier

Testing Data

Unseen Data

(Jeff, Professor, 4)

Tenured?

**Yes**

| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Tom | Assistant Prof | 2 | no |
| Merlisa | Associate Prof | 7 | no |
| George | Professor | 5 | yes |
| Joseph | Assistant Prof | 7 | yes |

# Attribute Selection Measures

- While implementing a Decision tree, the main issue arises <span style="color:red">that how to select the best attribute for the root node and for sub-nodes</span>. So, to solve such problems there is a technique which is called as **Attribute selection measure or ASM.**

- By this measurement, we can easily select the best attribute for the nodes of the tree. There are two popular techniques for ASM, which are:
  - **Information Gain**
  - **Gini Index**

$$Entropy\ (P) = -\sum_{i=1}^{n} p_i\ log_2\ (p_i)$$

# Information Gain and Gini Index in Decision Tree

$$Gini\ (P)\ =\ 1 - \sum_{i=1}^{n} (p_i)^2$$

# 1. Information Gain:

- Information gain is the **measurement of changes in entropy after the segmentation of a dataset based** on an attribute.

- It **calculates how much information** a feature **provides** us about a class.

**Entropy:** Entropy is a metric to measure the impurity in a given attribute. It specifies randomness in data. Entropy can be calculated as:

Entropy(s)= -P(yes)$\log_2$ P(yes)- P(no) $\log_2$ P(no)

**Where,**
- **S= Total number of samples**
- **P(yes)= probability of yes**
- **P(no)= probability of no**

$$Entropy = -p \log_2 p - q \log_2 q$$

$$Entropy = -0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$$

# Attribute Selection Measure: Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain

- Let $p_i$ be the probability that an arbitrary tuple in D belongs to class $C_i$, estimated by $|C_{i,D}|/|D|$

- Expected information (entropy) needed to classify a tuple in D:

$$Info(D) = -\sum_{i=1}^{m} p_i \log_2(p_i)$$
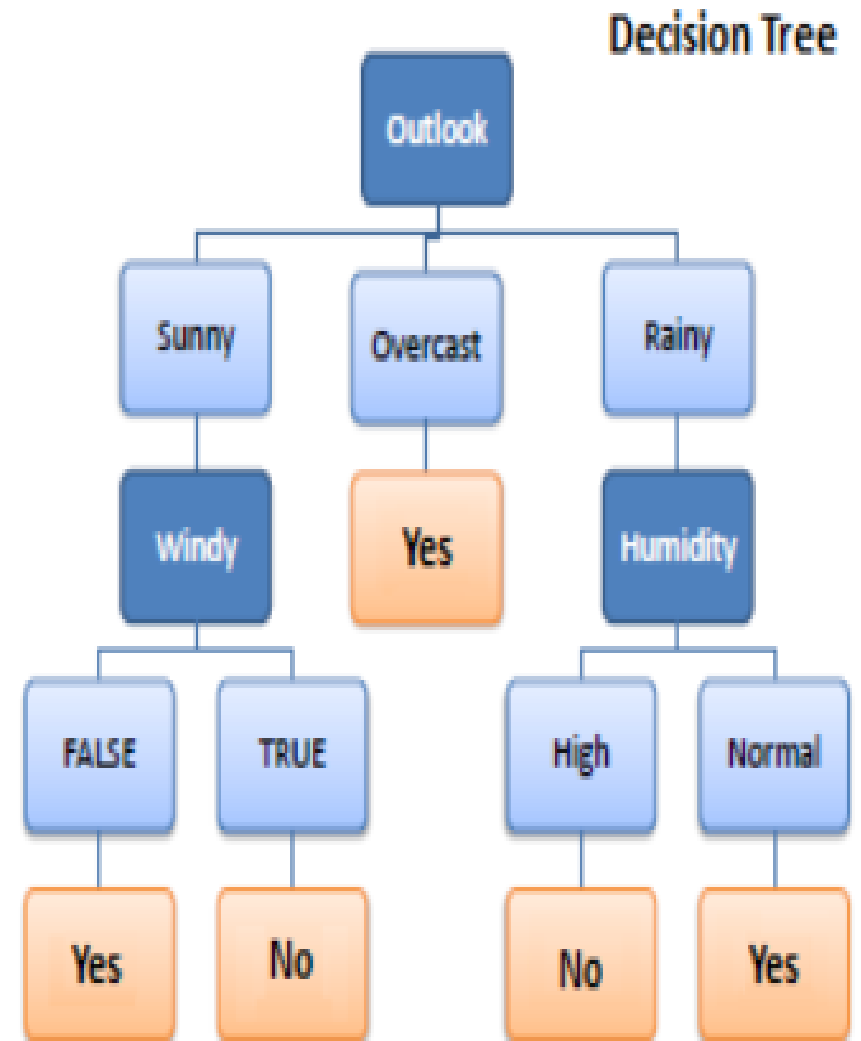
- Information needed (after using A to split D into v partitions) to classify D:

$$Info_A(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times Info(D_j)$$

- Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

Predictors     Target

| Outlook | Temp. | Humidity | Windy | Play Golf |
|---------|-------|----------|-------|-----------|
| Rainy | Hot | High | False | No |
| Rainy | Hot | High | True | No |
| Overcast | Hot | High | False | Yes |
| Sunny | Mild | High | False | Yes |
| Sunny | Cool | Normal | False | Yes |
| Sunny | Cool | Normal | True | No |
| Overcast | Cool | Normal | True | Yes |
| Rainy | Mild | High | False | No |
| Rainy | Cool | Normal | False | Yes |
| Sunny | Mild | Normal | False | Yes |
| Rainy | Mild | Normal | True | Yes |
| Overcast | Mild | High | True | Yes |
| Overcast | Hot | Normal | False | Yes |
| Sunny | Mild | High | True | No |

**Decision Tree**

$$\text{Entropy} = -p \log_2 p - q \log_2 q$$

$$\text{Entropy} = -0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$$

a) Entropy using the frequency table of one attribute:

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

| Play Golf | |
|-----------|-----------|
| Yes | No |
| 9 | 5 |

Entropy(PlayGolf) = Entropy (5,9)
      = Entropy (0.36, 0.64)
      = - (0.36 $\log_2$ 0.36) - (0.64 $\log_2$ 0.64)
      = 0.94

b) Entropy using the frequency table of two attributes:

$$E(T,X) = \sum_{c \in X} P(c)E(c)$$

| | | Play Golf | | |
|---|---|---|---|---|
| | | Yes | No | |
| | Sunny | 3 | 2 | 5 |
| Outlook | Overcast | 4 | 0 | 4 |
| | Rainy | 2 | 3 | 5 |
| | | | | 14 |

E(PlayGolf, Outlook) = P(Sunny)*E(3,2) + P(Overcast)*E(4,0) + P(Rainy)*E(2,3)

= (5/14)*0.971 + (4/14)*0.0 + (5/14)*0.971

= 0.693

| | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Outlook | Sunny | 3 | 2 |
| | Overcast | 4 | 0 |
| | Rainy | 2 | 3 |
| Gain = 0.247 | | | |

| | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Temp. | Hot | 2 | 2 |
| | Mild | 4 | 2 |
| | Cool | 3 | 1 |
| Gain = 0.029 | | | |

| | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Humidity | High | 3 | 4 |
| | Normal | 6 | 1 |
| Gain = 0.152 | | | |

| | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Windy | False | 6 | 2 |
| | True | 3 | 3 |
| Gain = 0.048 | | | |

$$Gain(T, X) = Entropy(T) - Entropy(T, X)$$

G(PlayGolf, Outlook) = E(PlayGolf) − E(PlayGolf, Outlook)

= 0.940 − 0.693 = 0.247

*Step 3*: Choose attribute with the largest information gain as the decision node, divide the dataset by its branches and repeat the same process on every branch.
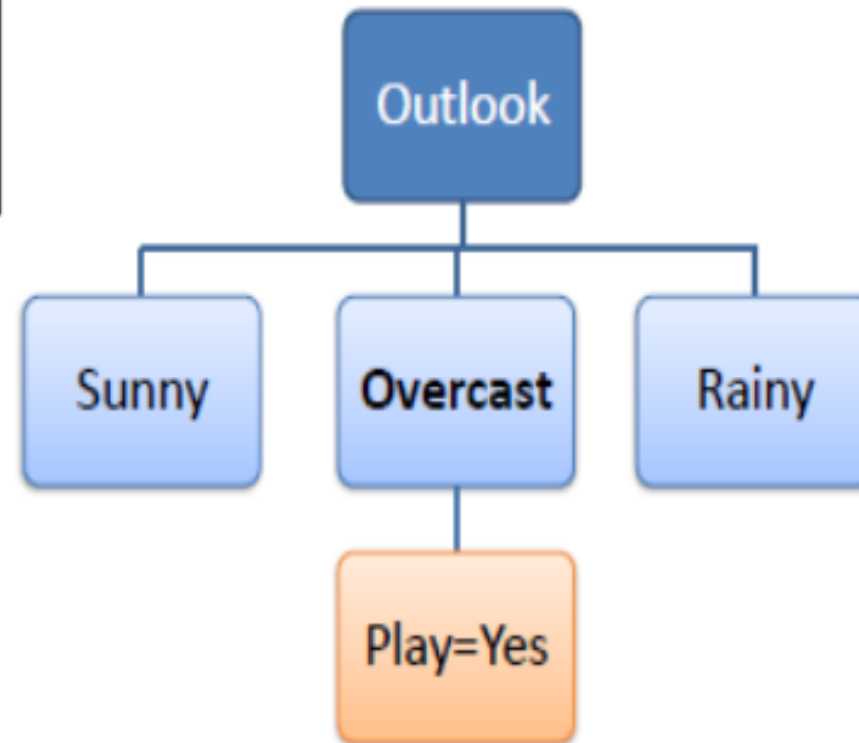
| | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Outlook | Sunny | 3 | 2 |
| | Overcast | 4 | 0 |
| | Rainy | 2 | 3 |
| Gain = 0.247 | | | |

**Sunny**

| Outlook | Temp | Humidity | Windy | Play Golf |
|---|---|---|---|---|
| Sunny | Mild | High | FALSE | Yes |
| Sunny | Cool | Normal | FALSE | Yes |
| Sunny | Cool | Normal | TRUE | No |
| Sunny | Mild | Normal | FALSE | Yes |
| Sunny | Mild | High | TRUE | No |

**Overcast**

| Overcast | Hot | High | FALSE | Yes |
|---|---|---|---|---|
| Overcast | Cool | Normal | TRUE | Yes |
| Overcast | Mild | High | TRUE | Yes |
| Overcast | Hot | Normal | FALSE | Yes |

**Rainy**

| Rainy | Hot | High | FALSE | No |
|---|---|---|---|---|
| Rainy | Hot | High | TRUE | No |
| Rainy | Mild | High | FALSE | No |
| Rainy | Cool | Normal | FALSE | Yes |
| Rainy | Mild | Normal | TRUE | Yes |

**Outlook**

*Step 4a*: A branch with entropy of 0 is a leaf node.

| Temp. | Humidity | Windy | Play Golf |
|-------|----------|-------|-----------|
| Hot | High | FALSE | Yes |
| Cool | Normal | TRUE | Yes |
| Mild | High | TRUE | Yes |
| Hot | Normal | FALSE | Yes |

*Step 4b*: A branch with entropy more than 0 needs further splitting.

| Temp | Humidity | Windy | Play Golf |
|------|----------|-------|-----------|
| Mild | High | FALSE | Yes |
| Cool | Normal | FALSE | Yes |
| Mild | Normal | FALSE | Yes |
| Cool | Normal | TRUE | No |
| Mild | High | TRUE | No |

*Step 5*: The ID3 algorithm is run recursively on the non-leaf branches, until all data is classified.

## Decision Tree to Decision Rules

A decision tree can easily be transformed to a set of rules by mapping from the root node to the leaf nodes one by one.

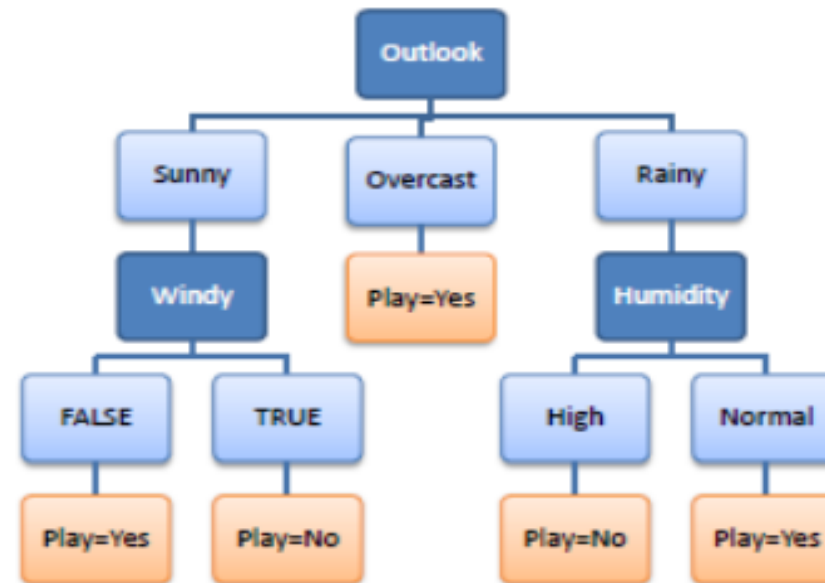**R₁**: IF (Outlook=Sunny) AND (Windy=FALSE) THEN Play=Yes

**R₂**: IF (Outlook=Sunny) AND (Windy=TRUE) THEN Play=No

**R₃**: IF (Outlook=Overcast) THEN Play=Yes

**R₄**: IF (Outlook=Rainy) AND (Humidity=High) THEN Play=No

**R₅**: IF (Outlook=Rain) AND (Humidity=Normal) THEN Play=Yes

# Homework

| ID | Fever | Cough | Breathing issues | Infected |
|----|-------|-------|------------------|----------|
| 1  | NO    | NO    | NO               | NO       |
| 2  | YES   | YES   | YES              | YES      |
| 3  | YES   | YES   | NO               | NO       |
| 4  | YES   | NO    | YES              | YES      |
| 5  | YES   | YES   | YES              | YES      |
| 6  | NO    | YES   | NO               | NO       |
| 7  | YES   | NO    | YES              | YES      |
| 8  | YES   | NO    | YES              | YES      |
| 9  | NO    | YES   | YES              | YES      |
| 10 | YES   | YES   | NO               | YES      |
| 11 | NO    | YES   | NO               | NO       |
| 12 | NO    | YES   | YES              | YES      |
| 13 | NO    | YES   | YES              | NO       |
| 14 | YES   | YES   | NO               | NO       |