



# CLASSIFICATION

T.R.Lekhaa  
AP/IT



# OUTLINE

- Introduction
- Applications
- Decision tree induction
- Algorithm
- Rule based classification



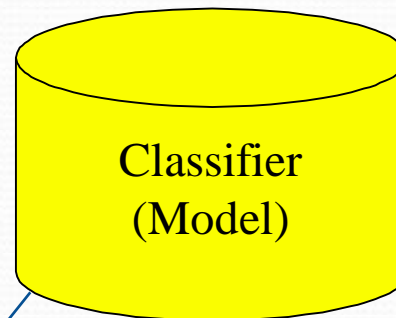
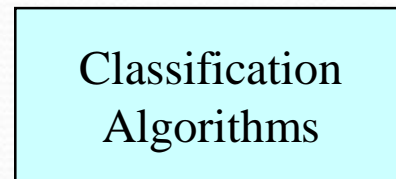
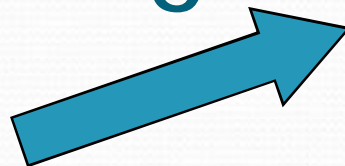
# Introduction

2-Steps:

- Learning
- Classification



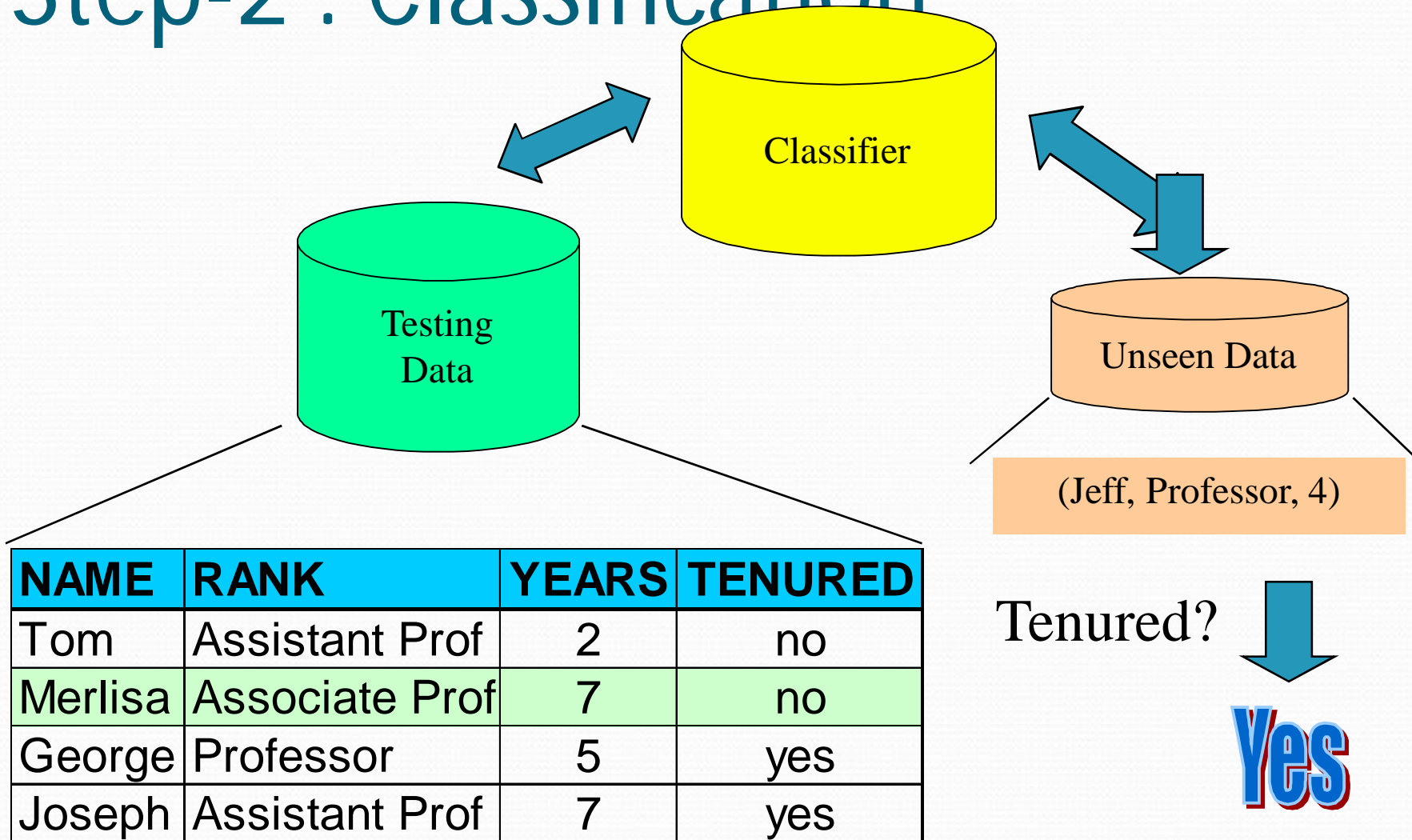
# Step-1 : Learning



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

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

# Step-2 : Classification



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



# Applications

- Credit card approval
- Target marketing
- Medical diagnosis



# Classification by decision tree

## Induction

Decision tree:

- Flowchart
- Internal node
- Branch
- Leaf node

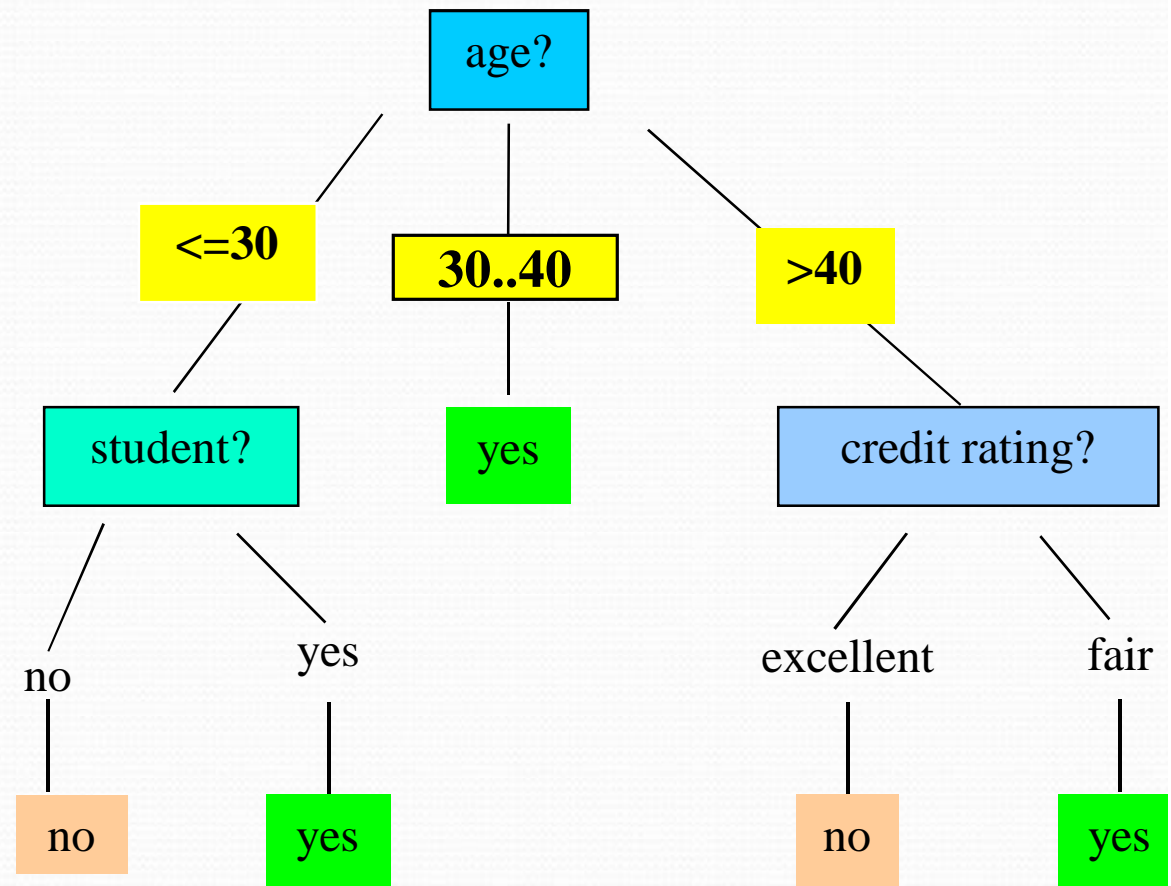
2-Phases:

- Tree construction
- Tree pruning



# Example

age	income	student	credit_rating
<=30	high	no	fair
<=30	high	no	excellent
31...40	high	no	fair
>40	medium	no	fair
>40	low	yes	fair
>40	low	yes	excellent
31...40	low	yes	excellent
<=30	medium	no	fair
<=30	low	yes	fair
>40	medium	yes	fair
<=30	medium	yes	excellent
31...40	medium	no	excellent
31...40	high	yes	fair
>40	medium	no	excellent







# Algorithm

Greedy algorithm:

- Top-Down approach
- Training data at root
- Attribute for measurement
  1. Information gain(ID3/C4.5)
  2. Gini index (IBM intelligent miner)



# Information gain(ID3/C4.5)

- All attributes are categorical
- Modified for continuous value attribute
- Select the attribute for high information gain
- 2 classes, P&N

Samples  $p$  element of P class &  $n$  element of N class

$$I(p,n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$



# Cont..

- Information gain measure used to select S, samples

- $I(s_1, s_2, \dots, s_m) = - \sum_{i=1}^m p_i \log_2 p_i$

Where  $p_i$ , probability &  $m$ , distinct value

Information encoded in bits  $p_i = \frac{s_i}{S}$

Entropy,  $E(A) = \sum_{i=1}^v I(s_{1j}, s_{2j}, \dots, s_{mj}) \frac{s_{1j}, s_{2j}, \dots, s_{mj}}{S}$

Subset,  $I(s_{1j}, s_{2j}, \dots, s_{mj}) = - \sum_{i=1}^m p_{ij} \log_2 p_{ij}$

$P_{ij} = \frac{s_{ij}}{|s_{ij}|}$



# Cont..

- $\text{Gain}(A) = I(s1, s2, sm) - E(A)$
  - $I(s1, s2) = -9/14 \log_2 (9/14) - 5/14 \log_2 (5/14)$   
 $= 0.940$
- Age  $\leq 30$   
Age 31...40  
Age  $> 40$



# Gini index(IBM intelligent miner)

- Attribute having continuous value
- Possible split values for each attribute
- Modified for categorical attributes
- Data set T contains n classes,

$$\text{Gini}(t) = 1 - \sum_{j=1}^n p_j^2 \quad \text{Where } p_j, \text{ relative frequency } j \text{ in } T$$

Dataset splits into 2 subset T1 & T2 with N1 & N2

$$\text{Gini}_{\text{split}}(T) = \frac{N_1}{N} \text{gini}(T1) + \frac{N_2}{N} \text{gini}(T2)$$

$$\begin{aligned} \text{Gini}(d) &= 1 - (9/14)^2 - (5/14)^2 \\ &= 0.459 \end{aligned}$$

$$\begin{aligned} \text{gini}_{\text{income}}(\text{low,medium})(D) &= 10/14 \text{gini}(D1) + 4/14 \text{gini}(D2) \\ &= 10/14(1 - (7/10)^2 - (3/10)^2) + 4/14(1 - (2/4)^2 - (2/4)^2) \\ &= 0.443 = \text{gini}_{\text{income}}(\text{high})(D) \end{aligned}$$



# Tree Pruning

- Reflected noise branches are removed
- 2 types
  1. Post pruning
  2. Pre pruning



# Rule based classification

- IF THEN rules are used
- Each attribute- conjunction
- Example:

IF age $\leq$ 30 & student=no THEN buys-computer=no

IF age $\leq$ 30 & student=Yes THEN buys-computer=Yes