



## Module 03 – Exercise Class

# DECISION TREE

# Classification & Regression

Nguyen Quoc Thai

# Objectives

## Decision Tree

- ❖ Introduction
- ❖ Regression & Classification Problem
- ❖ Terminology

## Decision Tree for Regression

- ❖ Variance
- ❖ Sum of Squared Errors (SSE)
- ❖ Sklearn Library

## Decision Tree for Classification

- ❖ Constructing Decision Tree: ID3
- ❖ Gini Impurity
- ❖ Entropy
- ❖ Information Gain
- ❖ Sklearn library

## Improved Decision Tree

- ❖ Overfitting
- ❖ Stopping Early
- ❖ Post-Pruning (Reduced-Error & Rule Post-Pruning)
- ❖ Missing Attribute Values

## SECTION 1

## Decision Tree

## SECTION 2

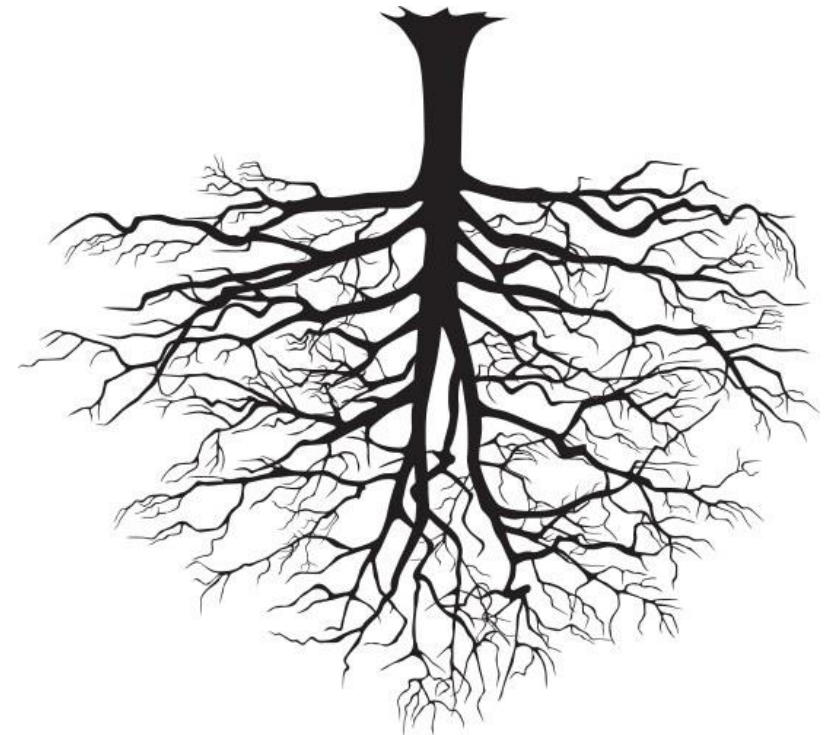
## DT for Classification

## SECTION 3

## DT for Regression

## SECTION 4

## Improved DT

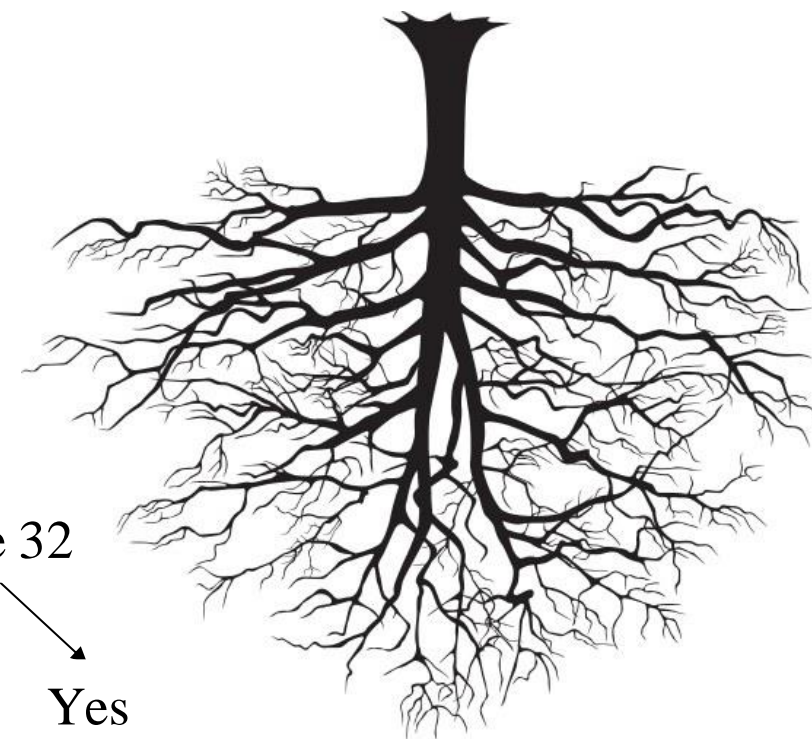
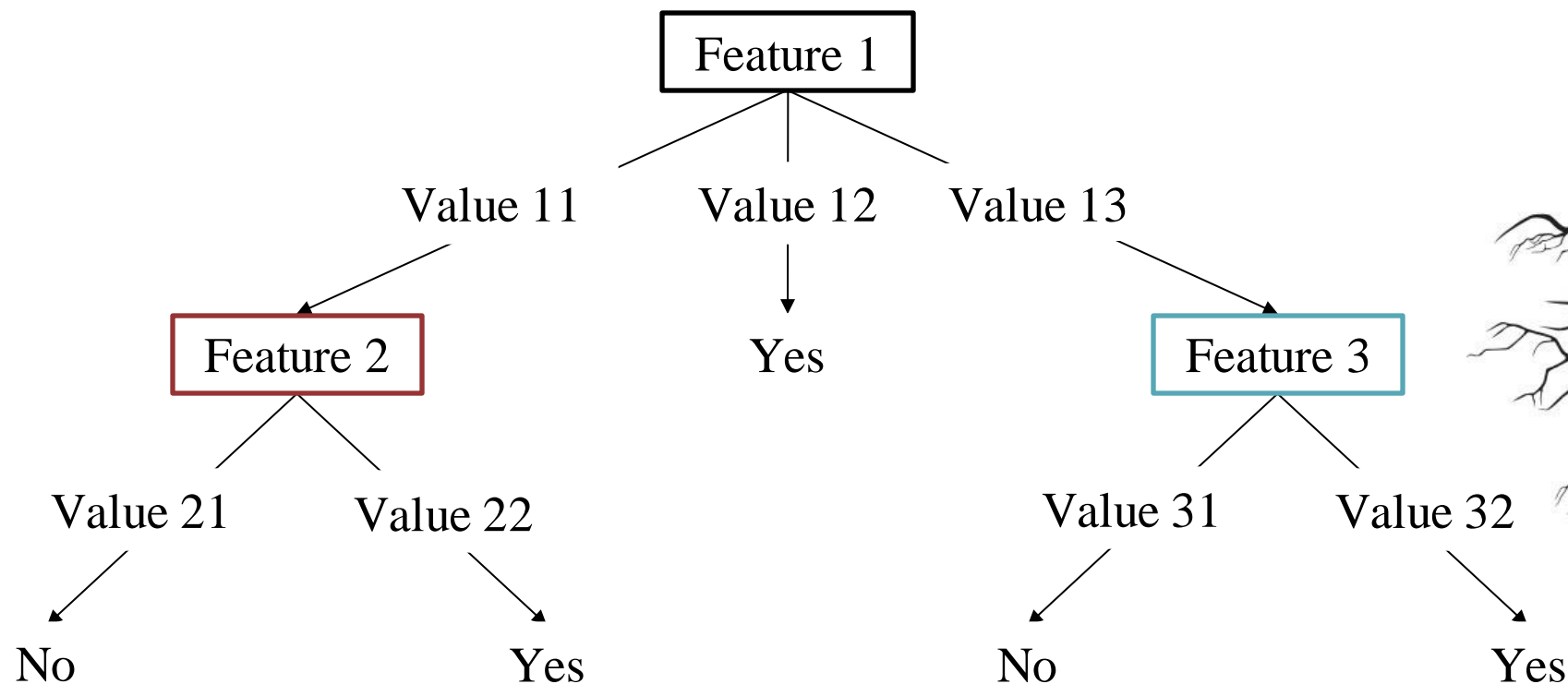


# Decision Tree



## Introduction

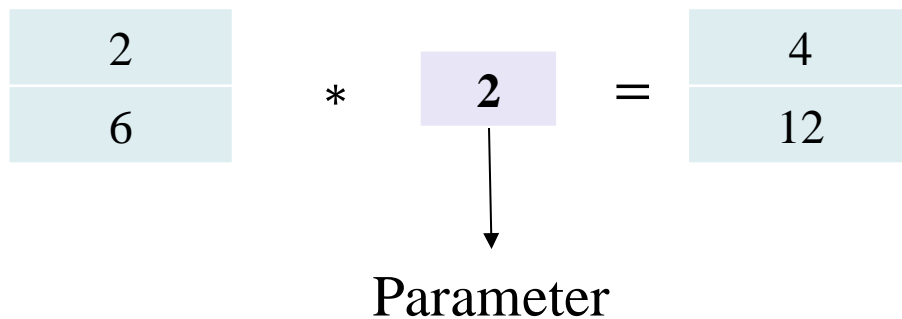
A Decision Tree is a non-parametric Supervised Machine Learning algorithm  
Used for both regression and classification problem



## ! Non-parametric Supervised Machine Learning

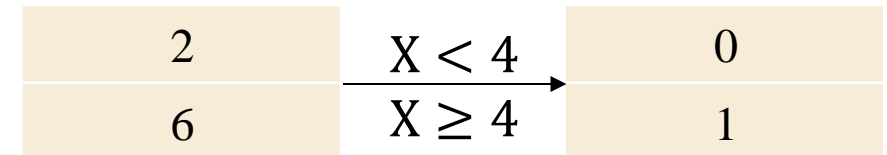
### Parametric SML

X	Y
2	4
6	12



### Non-parametric SML

X	Y
2	0
6	1





## Regression & Classification Problem

### Regression

- Predict a continuous value based on the input variables



What will be the temperature tomorrow?



### Classification

- Classify input variables to identify discrete output variables (labels, categories)



Will it be hot or cold tomorrow?



## ! Regression & Classification Problem

### Classification

- Classify input variables to identify discrete output variables (labels, categories)

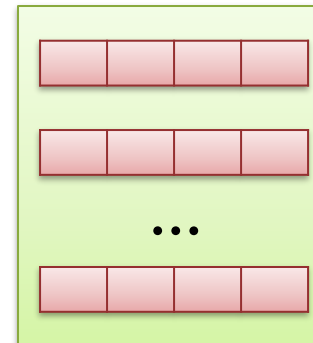


Will it be hot or cold tomorrow?



### Training Data

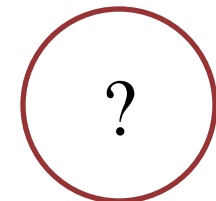
Feature



Label



### Test Data



## ! Regression & Classification Problem

### Regression

- Predict a continuous value based on the input variables

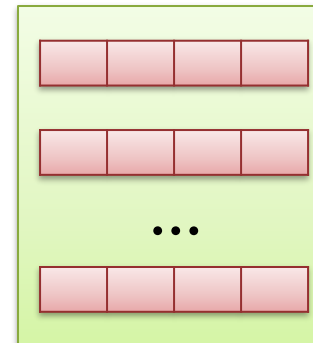


What will be the temperature tomorrow?



### Training Data

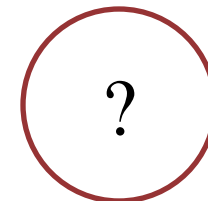
Feature



Continuous Label



### Test Data



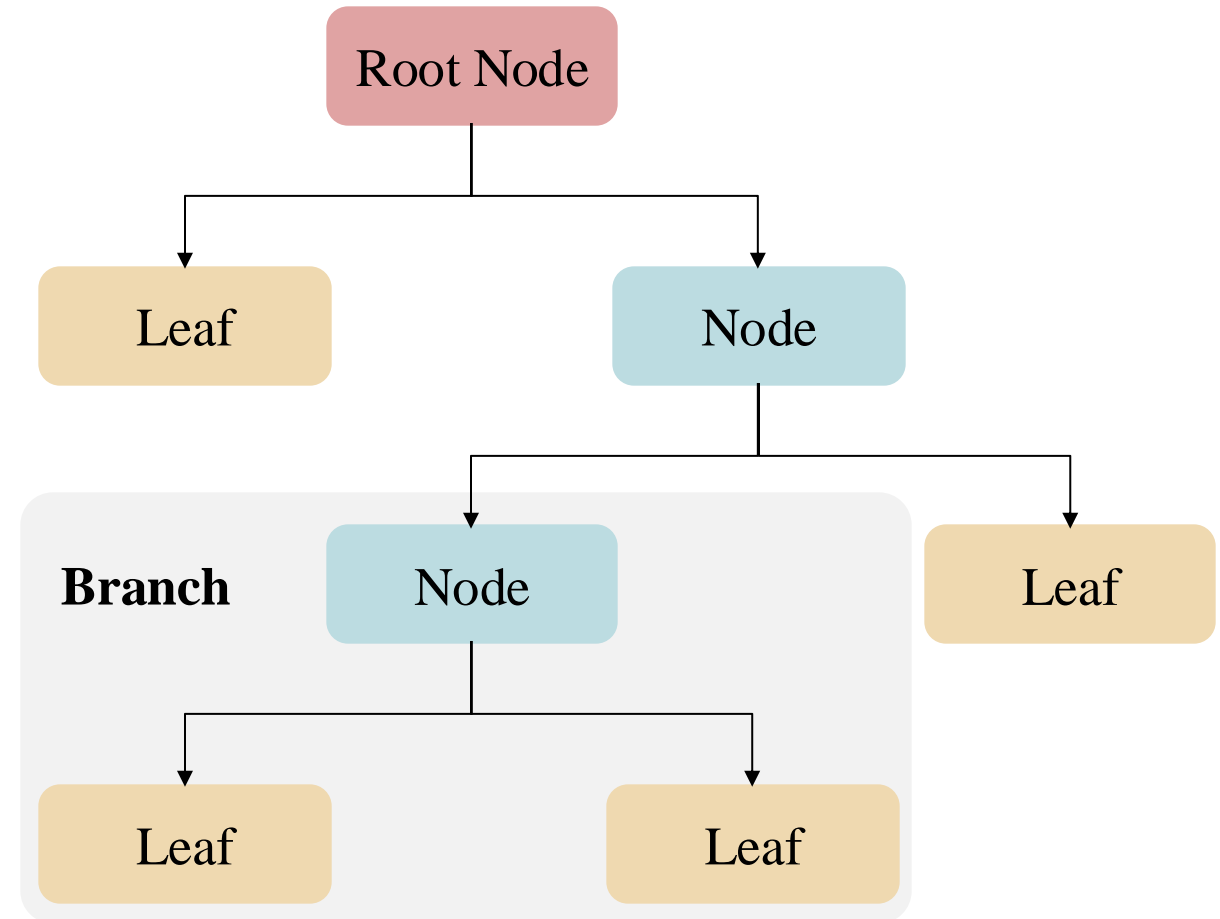


# Decision Tree



## Terminology

- ❖ **Root Node:** the top-level node
- ❖ **Node:** internal node or decision node
- ❖ **Parent Node:** a node that precedes a (child) node
- ❖ **Leaf:** terminal node – a node at the end of a branch – represents outcome of the tree (label or numerical value)
- ❖ **Branches:** a subset of a tree, starting at an (internal) node until the leaves

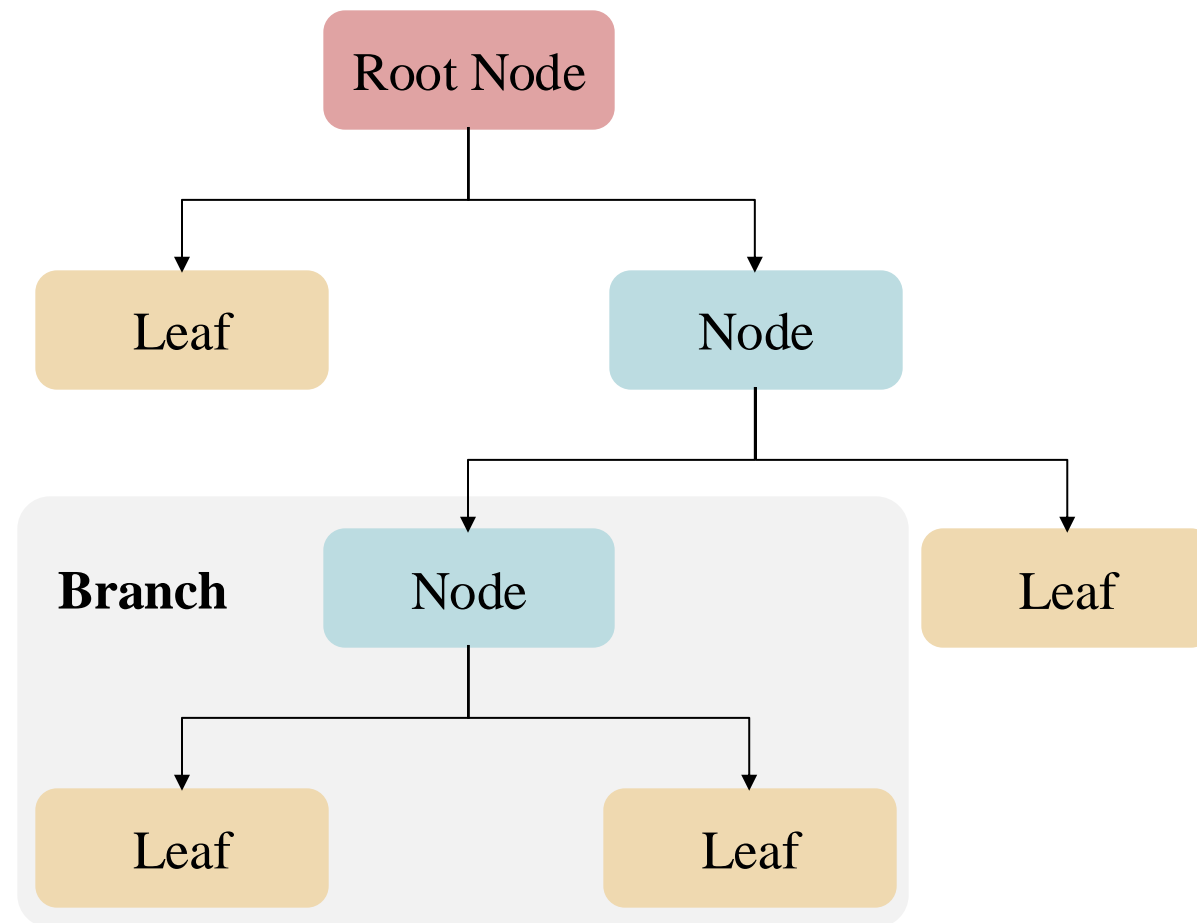


# Decision Tree



## Terminology

- ❖ **Splitting:** divide a node into two child nodes depending on a criterion and a selected feature
- ❖ **Pruning:** removing a branch from a tree
- ❖ Described by **IF-ELSE** Sets



## SECTION 1

### Decision Tree

## SECTION 2

### DT for Classification

## SECTION 3

### DT for Regression

## SECTION 4

### Improved DT

Dataset D (n sample)



Dataset D<sub>1</sub>  
(n<sub>1</sub> sample)

Dataset D<sub>2</sub>  
(n<sub>2</sub> sample)

Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification



## Decision Tree



### Example Dataset

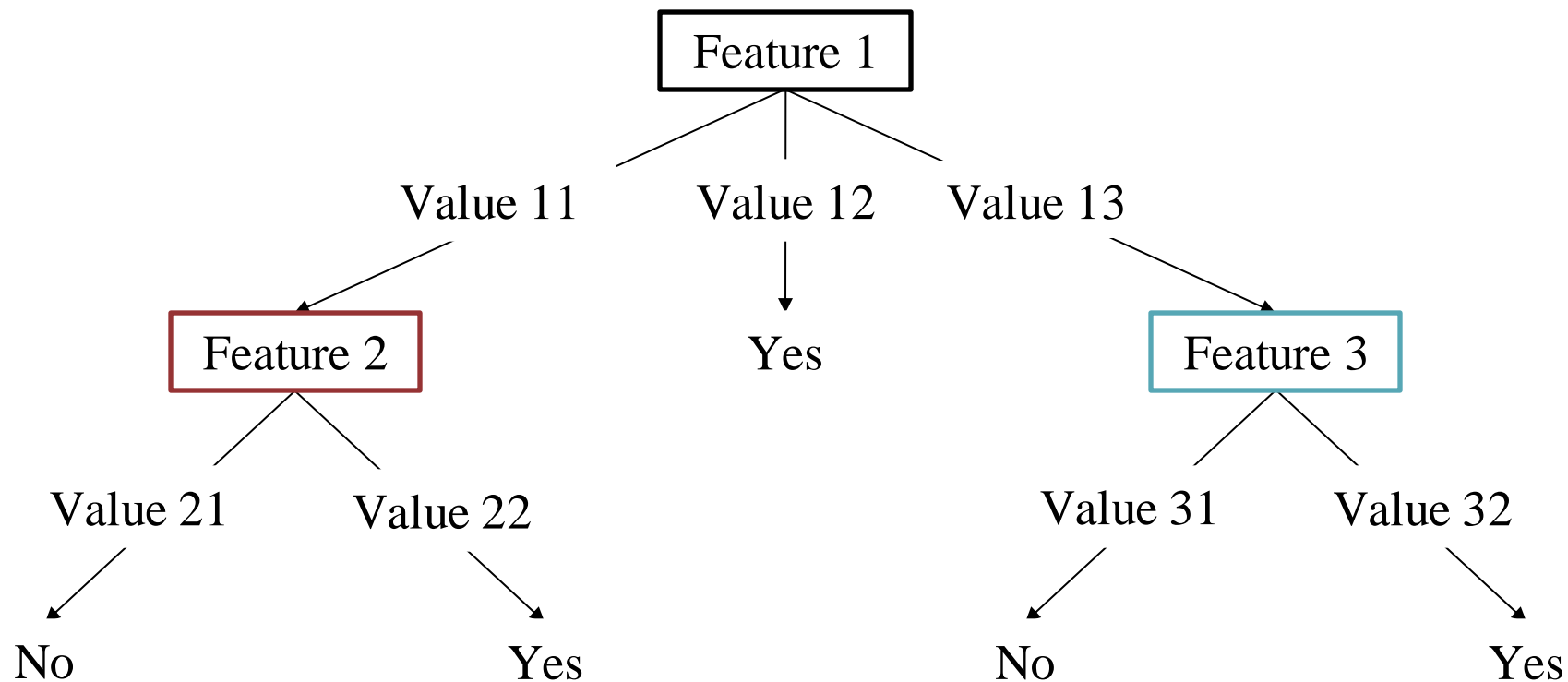
Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification



## Decision Tree

- Main step: splitting the data according to a splitting criterion



# DT for Classification

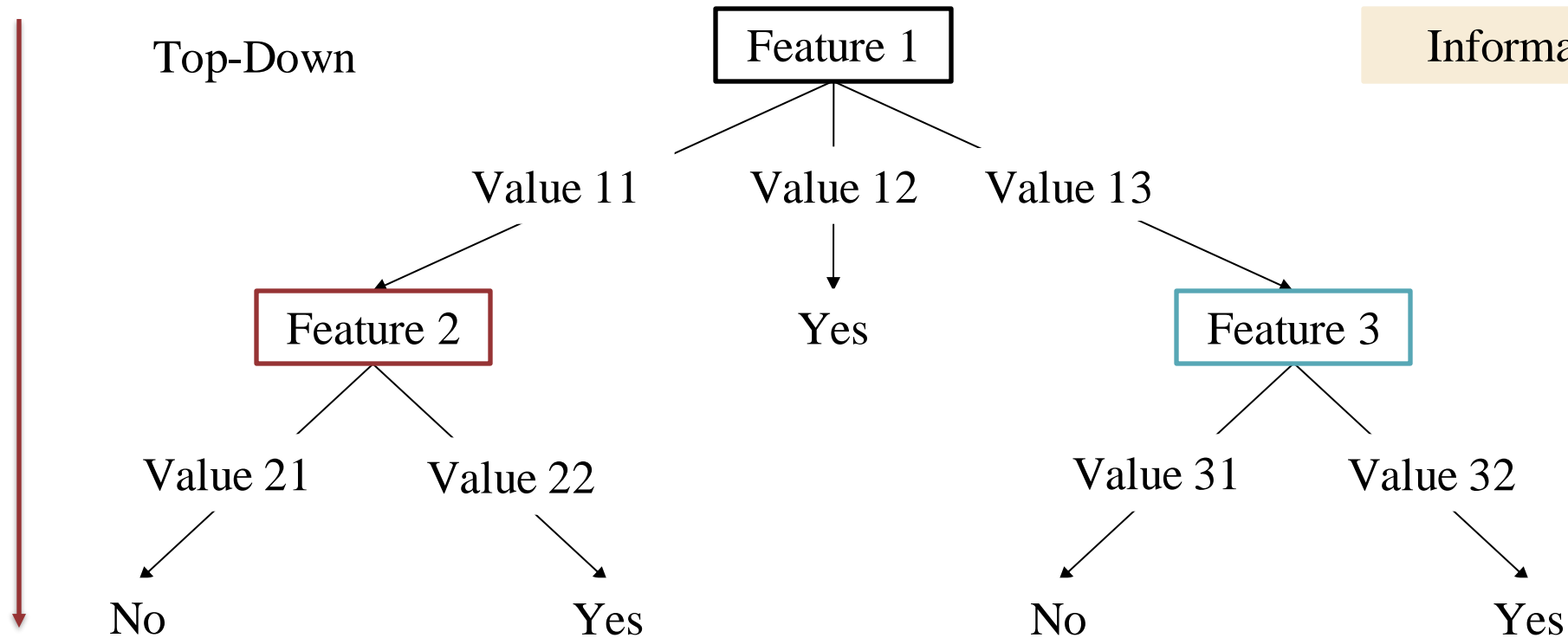


## Constructing Decision Tree: ID3

### ➤ Iterative Dichotomiser 3

Gini Impurity

Information Gain



# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

$p_j$  the probability that a random drawn sample from this node belongs to class  $j$  and  $c$  the number of classes

1. For each possible split, create child nodes and calculate the Gini Impurity of each child node.
2. Calculate the Gini Impurity of the split as the weighted average Gini Impurity of child nodes.
3. Select the split with the lowest Gini Impurity.

Dataset D (n sample)



Dataset D<sub>1</sub>  
(n<sub>1</sub> sample)

Dataset D<sub>2</sub>  
(n<sub>2</sub> sample)

Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

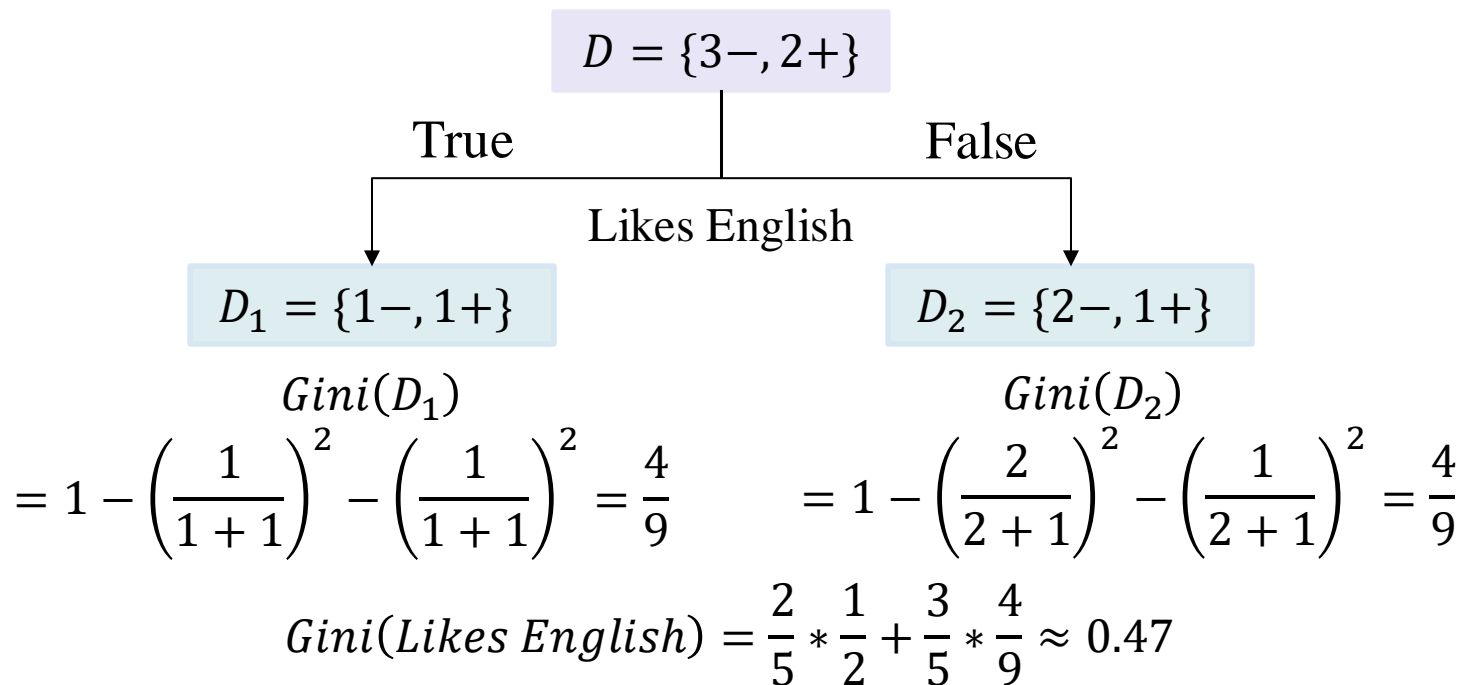
# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$



Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0



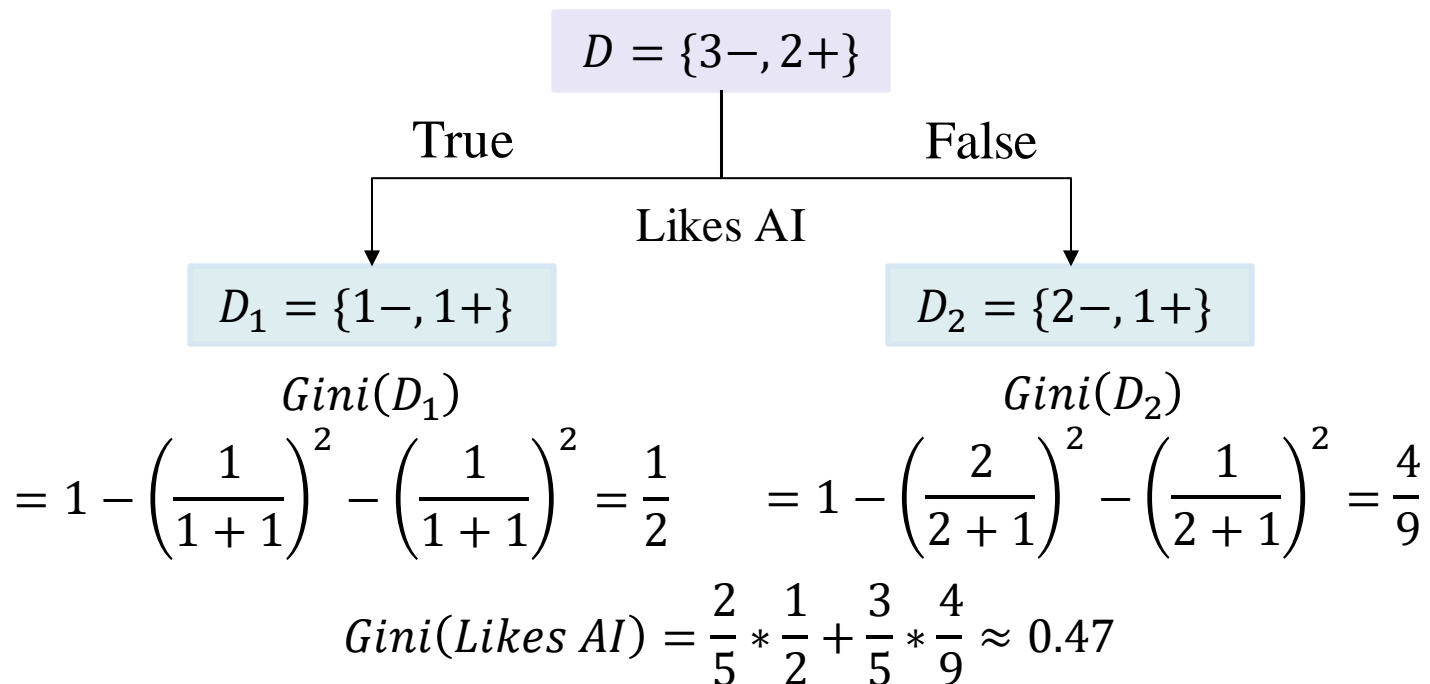
# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$



Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

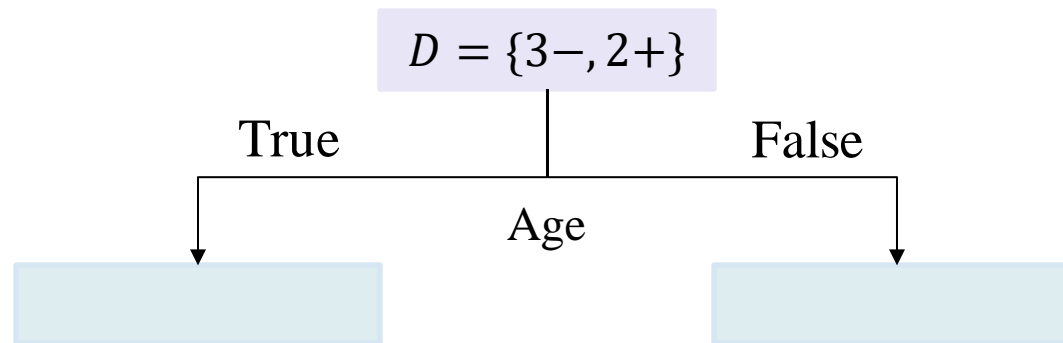
# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$



Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

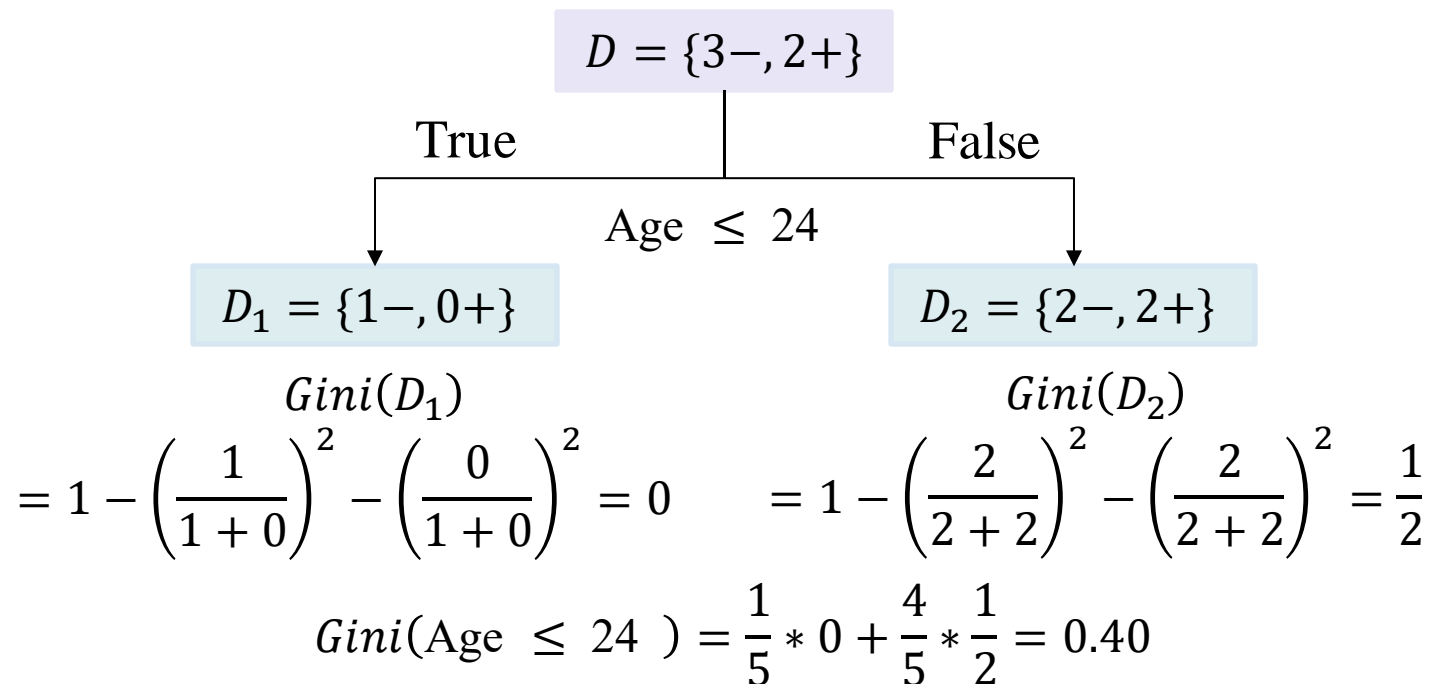
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



	Age	Likes English	Likes AI	Raise Salary
	23	0	0	<b>0</b>
<b>24</b>	25	1	1	<b>0</b>
<b>26</b>	27	1	0	<b>1</b>
<b>28</b>	29	0	1	<b>1</b>
<b>29</b>	29	0	0	<b>0</b>

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

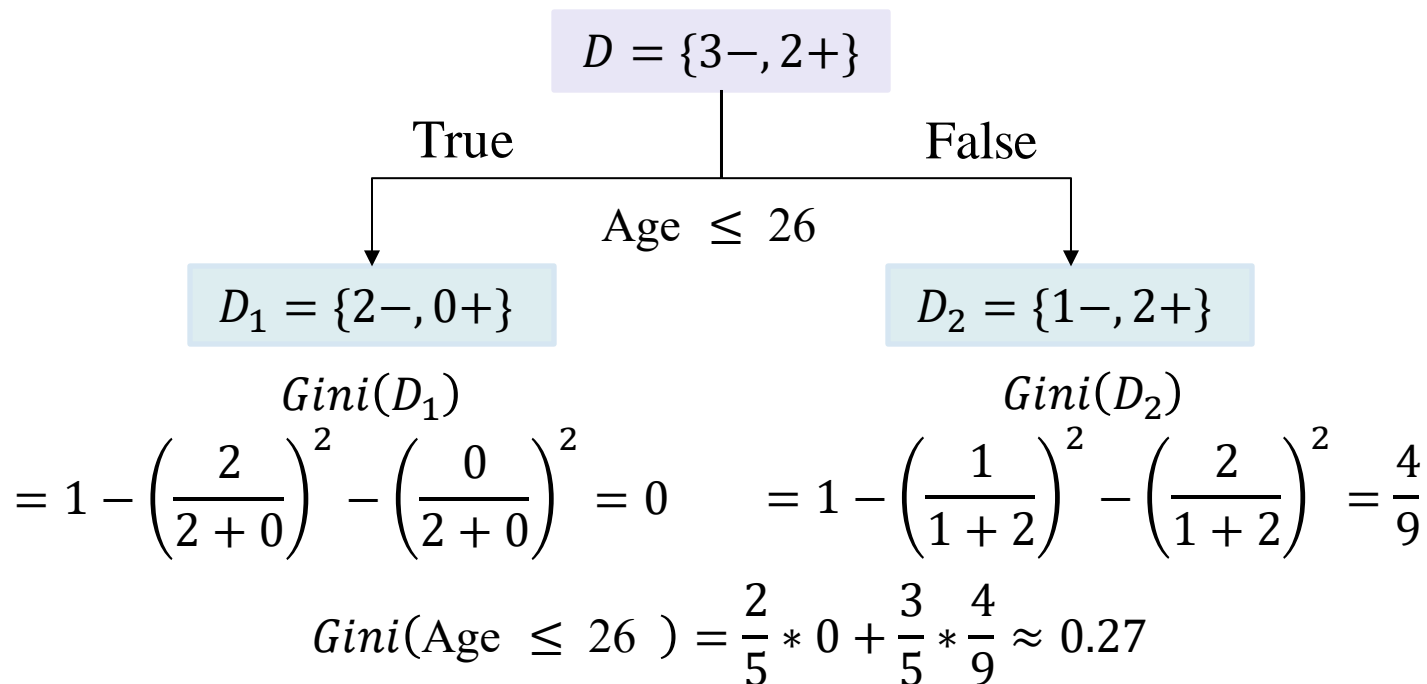
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	25	1	1	0
26	27	1	0	1
	29	0	1	1
28	29	0	0	0
	29	0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

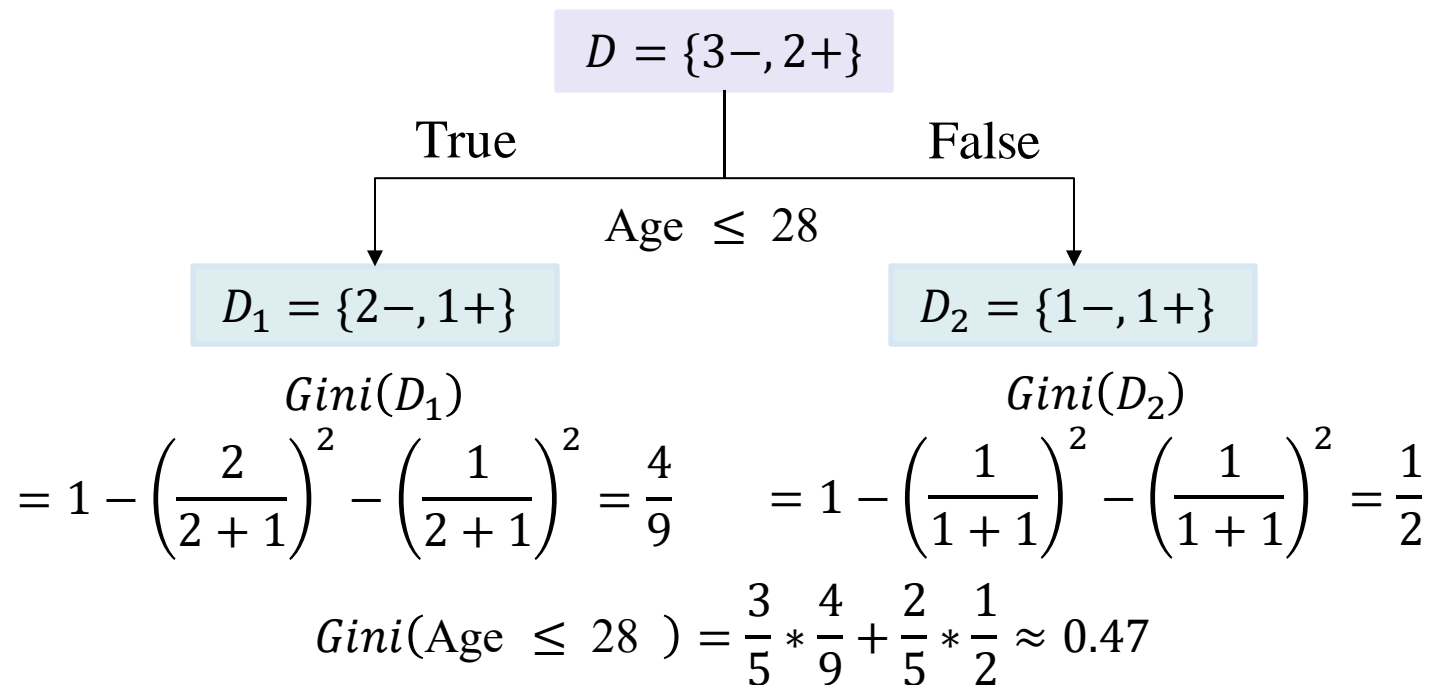
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	26	25	1	1
28		27	1	0
	29	29	0	1
29		0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

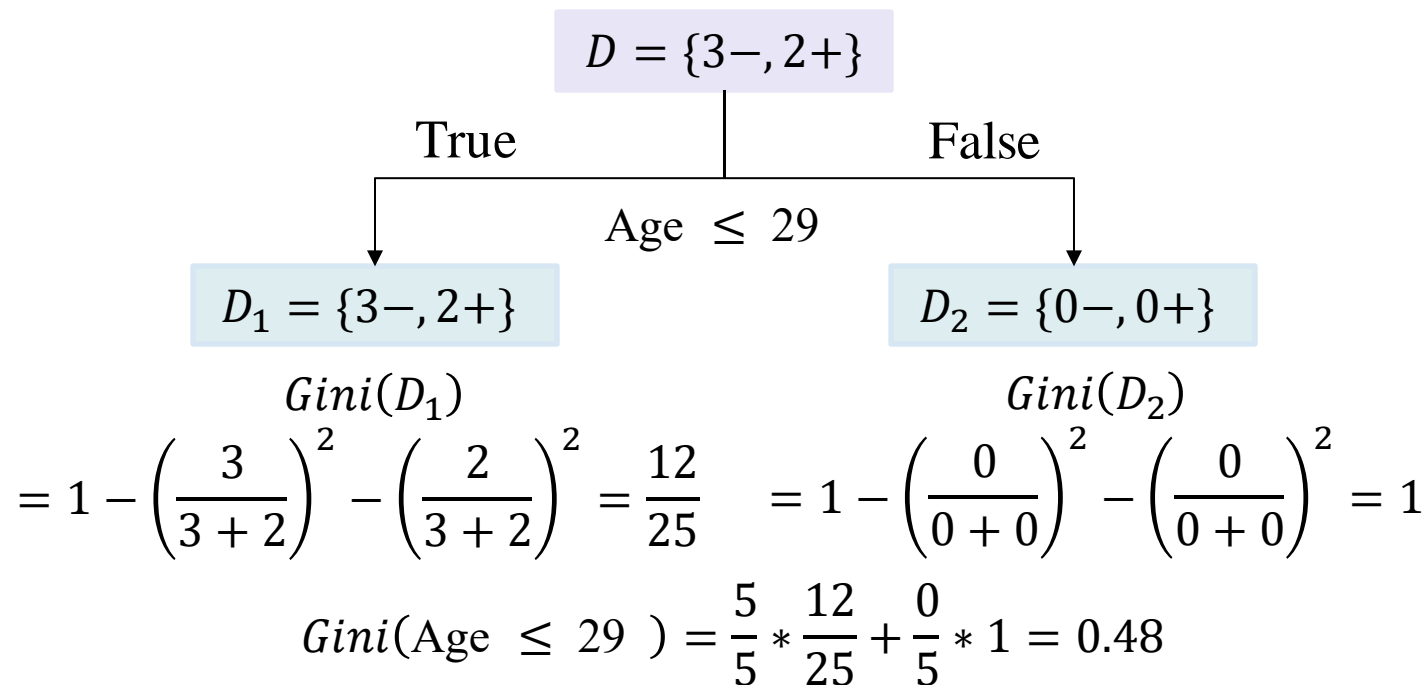
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
26	25	1	1	0
28	27	1	0	1
29	29	0	1	1
29	29	0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

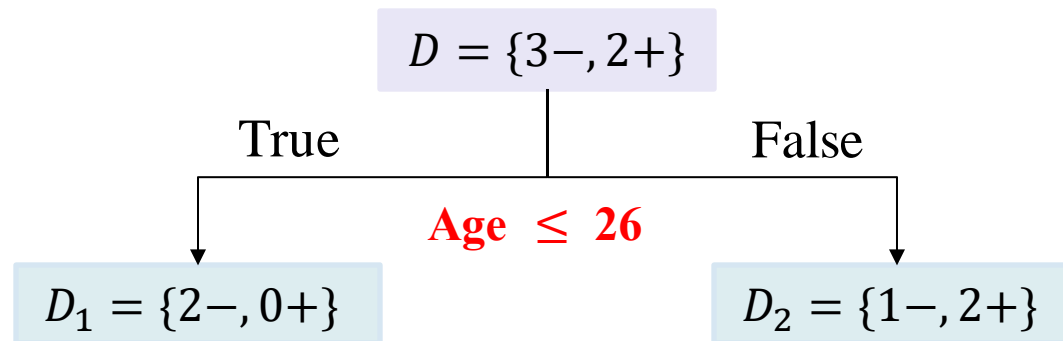
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



$$Gini(Likes\ English) \approx 0.47$$

$$Gini(Likes\ AI) \approx 0.47$$

$$Gini(Age \leq 24) = 0.40$$

$$Gini(Age \leq 26) = 0.27$$

$$Gini(Age \leq 28) = 0.47$$

$$Gini(Age \leq 29) = 0.48$$

	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	25	1	1	0
26	27	1	0	1
	29	0	1	1
28	29	0	0	0
	29	0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

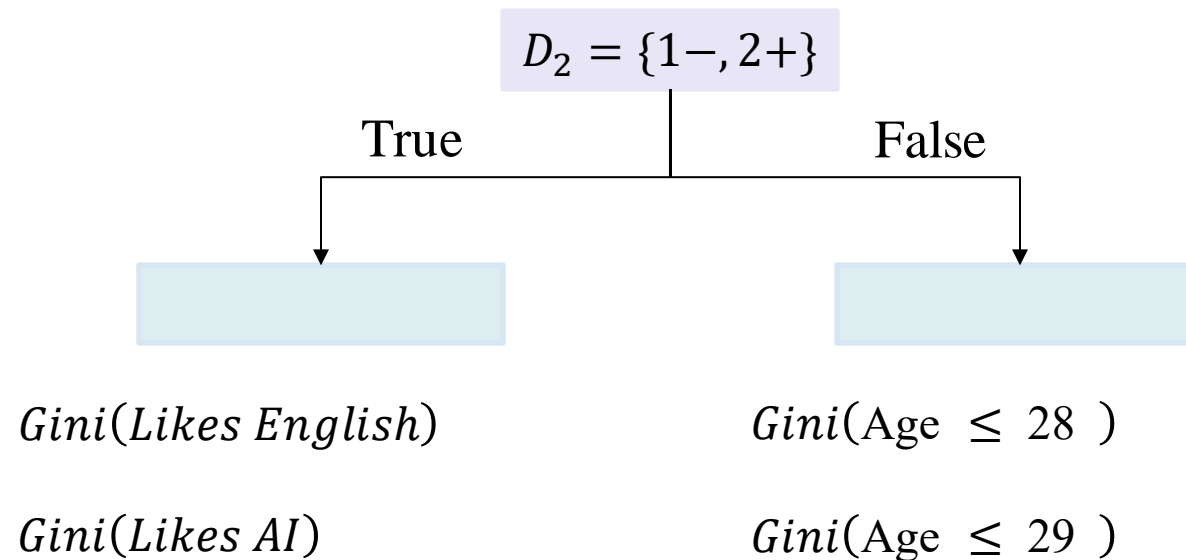
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	25	1	1	0
26	27	1	0	1
	29	0	1	1
28	29	0	0	0
	29	0	0	0



# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

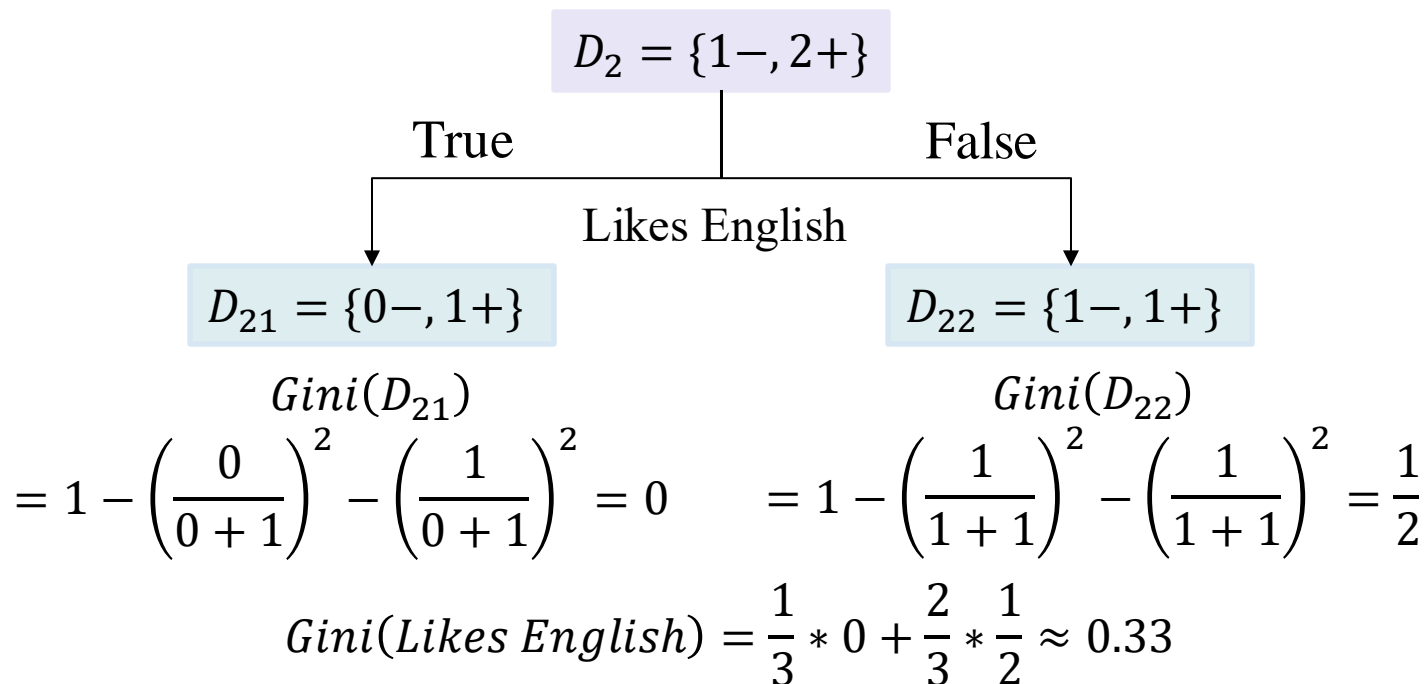
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	26	25	1	1
28		27	1	0
	29	29	0	1
29		0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

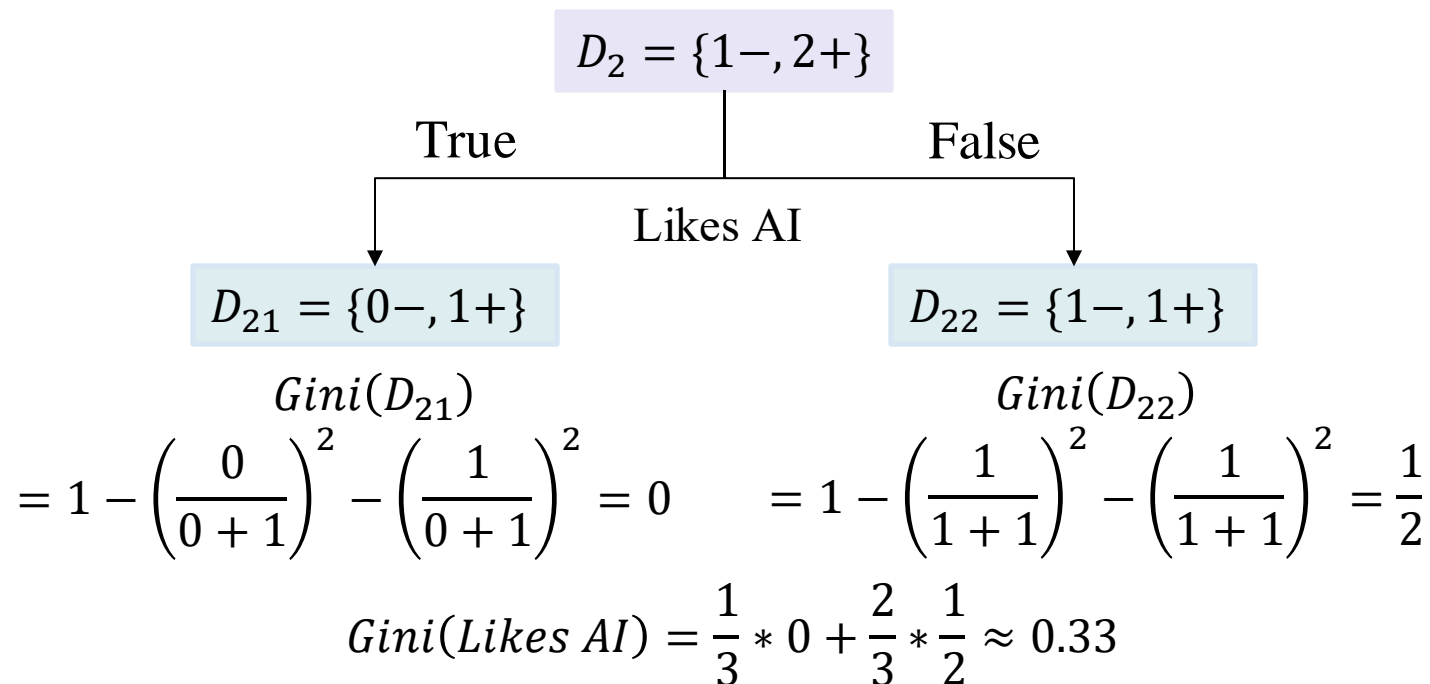
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	25	1	1	0
26	27	1	0	1
28	29	0	1	1
29	29	0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

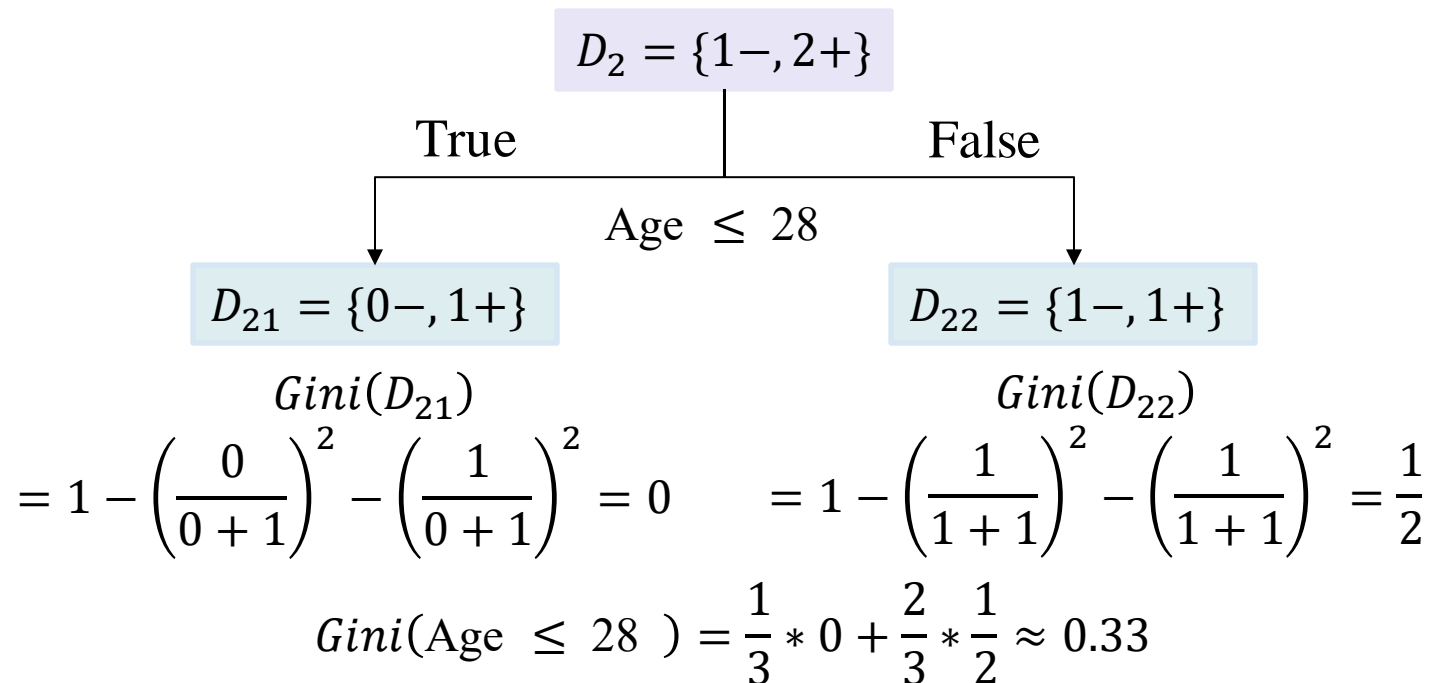
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	25	1	1	0
26	27	1	0	1
28		0	1	1
29	29	0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

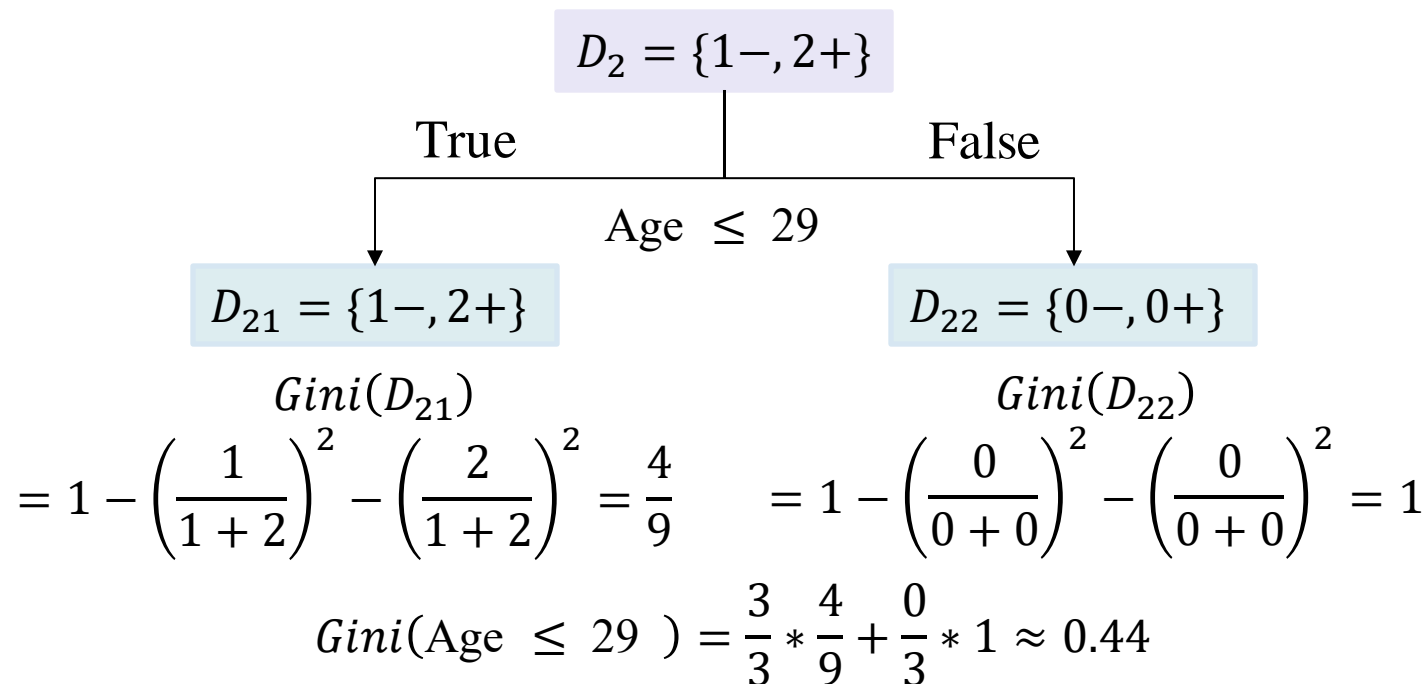
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	26	25	1	1
28		27	1	0
	29	29	0	1
29		0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

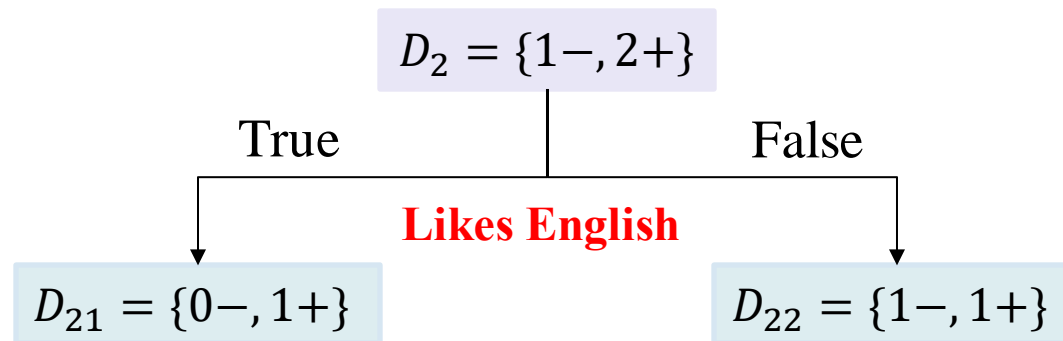
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



$$Gini(Likes\ English) \approx 0.33$$

$$Gini(Likes\ AI) \approx 0.33$$

$$Gini(Age \leq 28) \approx 0.33$$

$$Gini(Age \leq 29) \approx 0.44$$

	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	25	1	1	0
26	27	1	0	1
	29	0	1	1
28	29	0	0	0
	29	0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

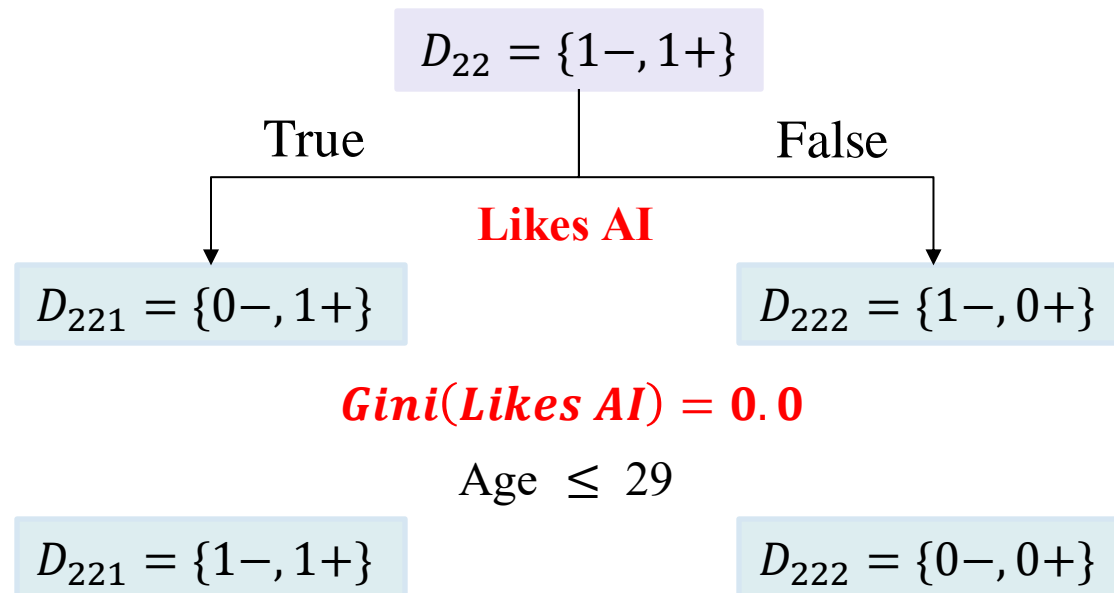
$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini



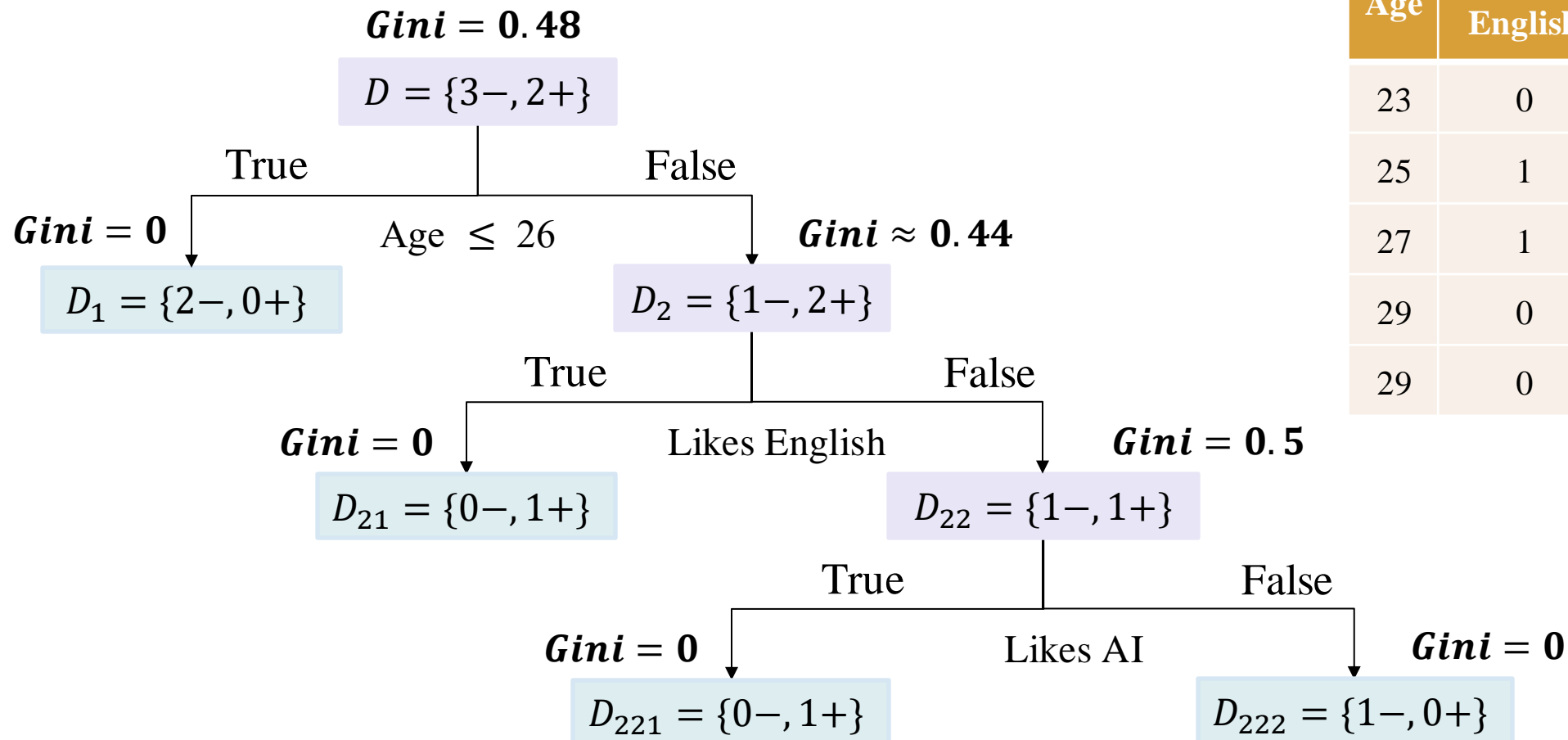
$$Gini(\text{Age} \leq 29) = 0.5$$

	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
	25	1	1	0
26	27	1	0	1
	28	0	1	1
29	29	0	0	0
	29	0	0	0

# DT for Classification



## Gini Impurity



Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification



## Gini Impurity

$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

$p_j$  the probability that a random drawn sample from this node belongs to class  $j$  and  $c$  the number of classes

1. For each possible split, create child nodes and calculate the Gini Impurity of each child node.
2. Calculate the Gini Impurity of the split as the weighted average Gini Impurity of child nodes.
3. Select the split with the lowest Gini Impurity.

Dataset D (n sample)



Dataset D<sub>1</sub>  
(n<sub>1</sub> sample)

Dataset D<sub>2</sub>  
(n<sub>2</sub> sample)

	Age	Likes English	Likes AI	Raise Salary
24	23	0	0	0
26	25	1	1	0
27	27	1	0	1
29	29	0	1	1
—	29	0	0	0



# DT for Classification



## Infomation Gain

$$Gain(D) = 1 - Entropy(D)$$

$$Entropy(D) = \frac{n_1}{n} Entropy(D_1) + \frac{n_2}{n} Entropy(D_2)$$

$$Entropy(D_i) = -\sum_{j=1}^c p_j \log_2 p_j$$

$p_j$  the probability that a random drawn sample from this node belongs to class  $j$  and  $c$  the number of classes

1. For each possible split, create child nodes and calculate the Entropy of each child node.
2. Calculate the Entropy of the split as the weighted average Entropy of child nodes.
3. Select the split with the lowest Entropy or highest Information Gain, respectively.

Dataset D (n sample)



Dataset D<sub>1</sub>  
(n<sub>1</sub> sample)

Dataset D<sub>2</sub>  
(n<sub>2</sub> sample)

Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification



## Information Gain

$$Gain(D) = 1 - Entropy(D)$$

$$Entropy(D) = \frac{n_1}{n} Entropy(D_1) + \frac{n_2}{n} Entropy(D_2)$$

$$Entropy(D_i) = - \sum_{j=1}^c p_j \log_2 p_j$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Entropy / Gain

$$D = \{3-, 2+\}$$

$$Entropy(D) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \approx 0.971$$

Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification

## ! Infomation Gain

$$Gain(D) = 1 - Entropy(D)$$

$$Entropy(D) = \frac{n_1}{n} Entropy(D_1) + \frac{n_2}{n} Entropy(D_2)$$

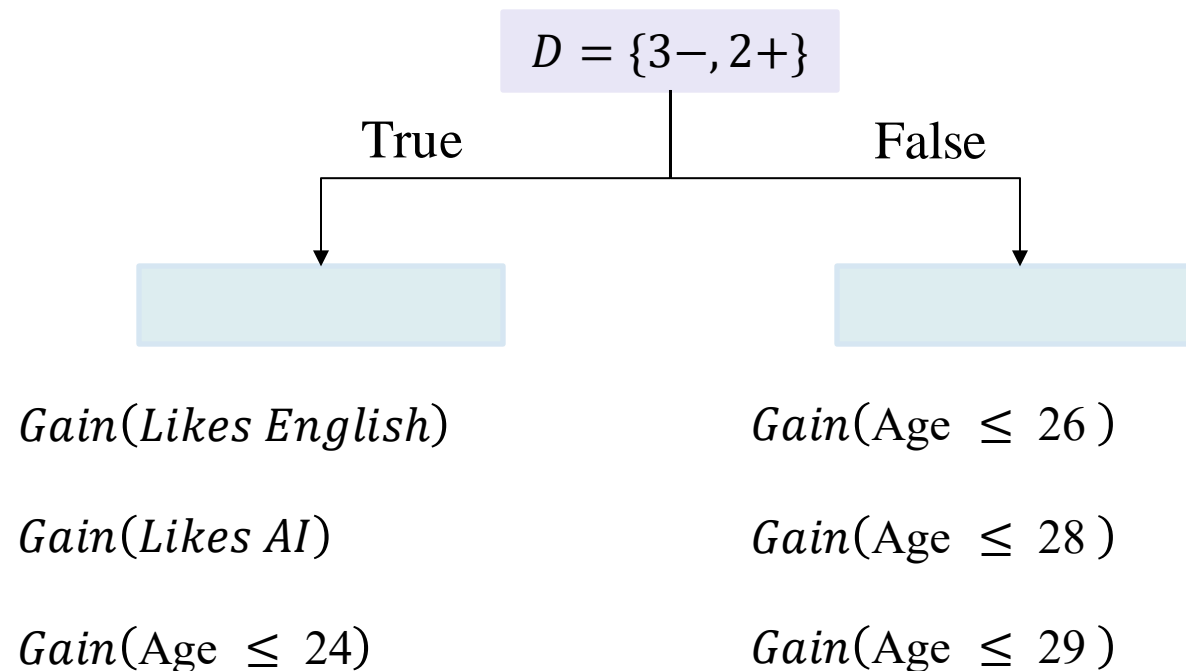
$$Entropy(D_i) = - \sum_{j=1}^c p_j \log_2 p_j$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Entropy / Gain



Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification

## ! Infomation Gain

$$Gain(D) = 1 - Entropy(D)$$

$$Entropy(D) = \frac{n_1}{n} Entropy(D_1) + \frac{n_2}{n} Entropy(D_2)$$

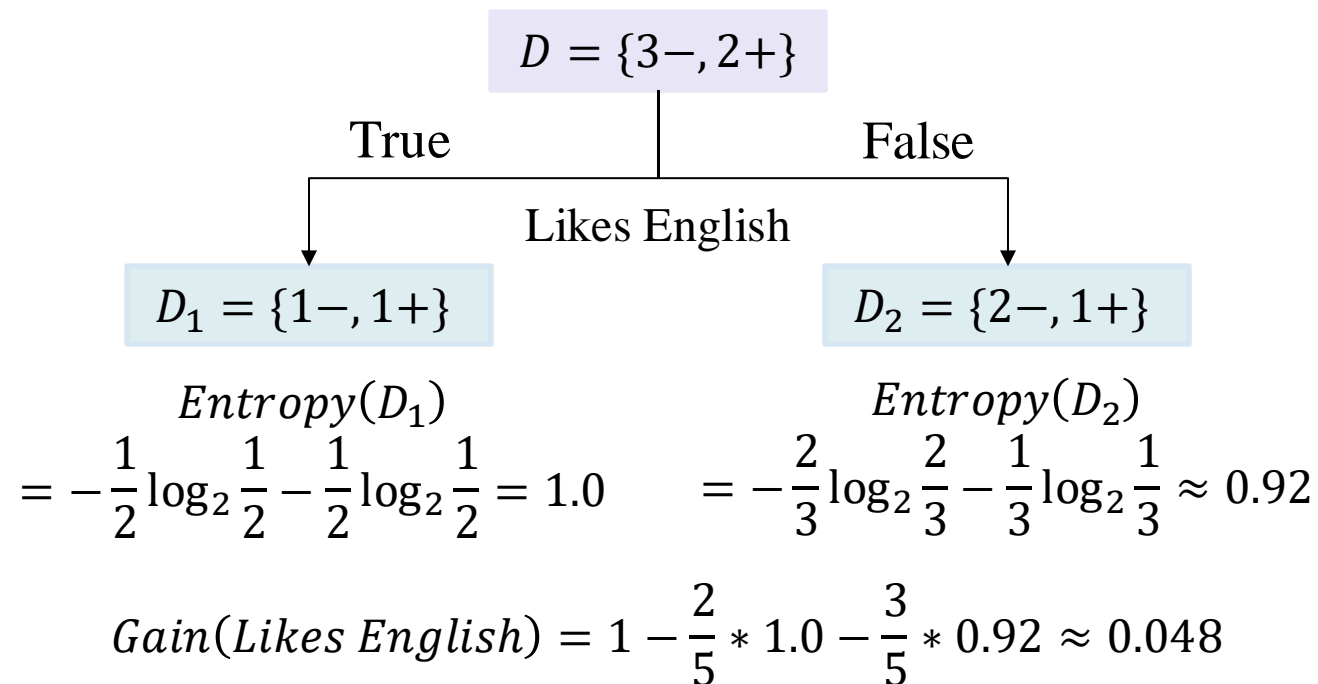
$$Entropy(D_i) = - \sum_{j=1}^c p_j \log_2 p_j$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Entropy / Gain



Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification

## ! Infomation Gain

$$Gain(D) = 1 - Entropy(D)$$

$$Entropy(D) = \frac{n_1}{n} Entropy(D_1) + \frac{n_2}{n} Entropy(D_2)$$

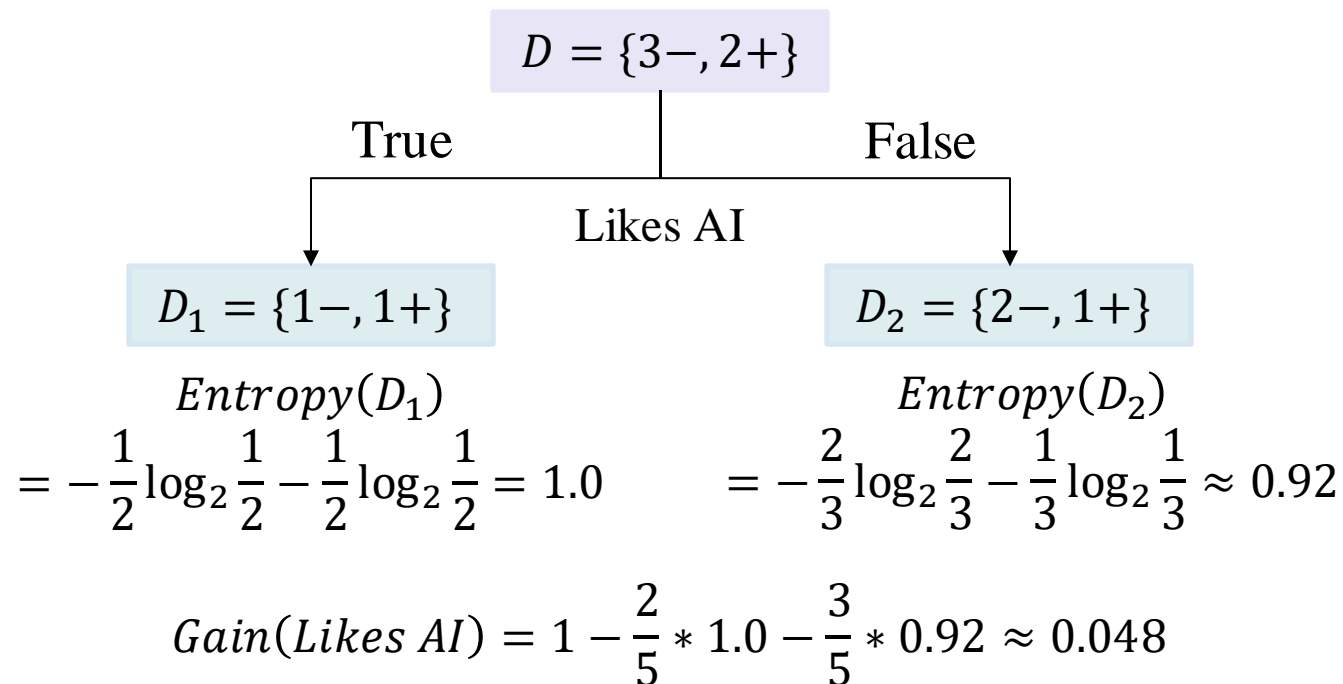
$$Entropy(D_i) = - \sum_{j=1}^c p_j \log_2 p_j$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Entropy / Gain



Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification

## ! Infomation Gain

$$Gain(D) = 1 - Entropy(D)$$

$$Entropy(D) = \frac{n_1}{n} Entropy(D_1) + \frac{n_2}{n} Entropy(D_2)$$

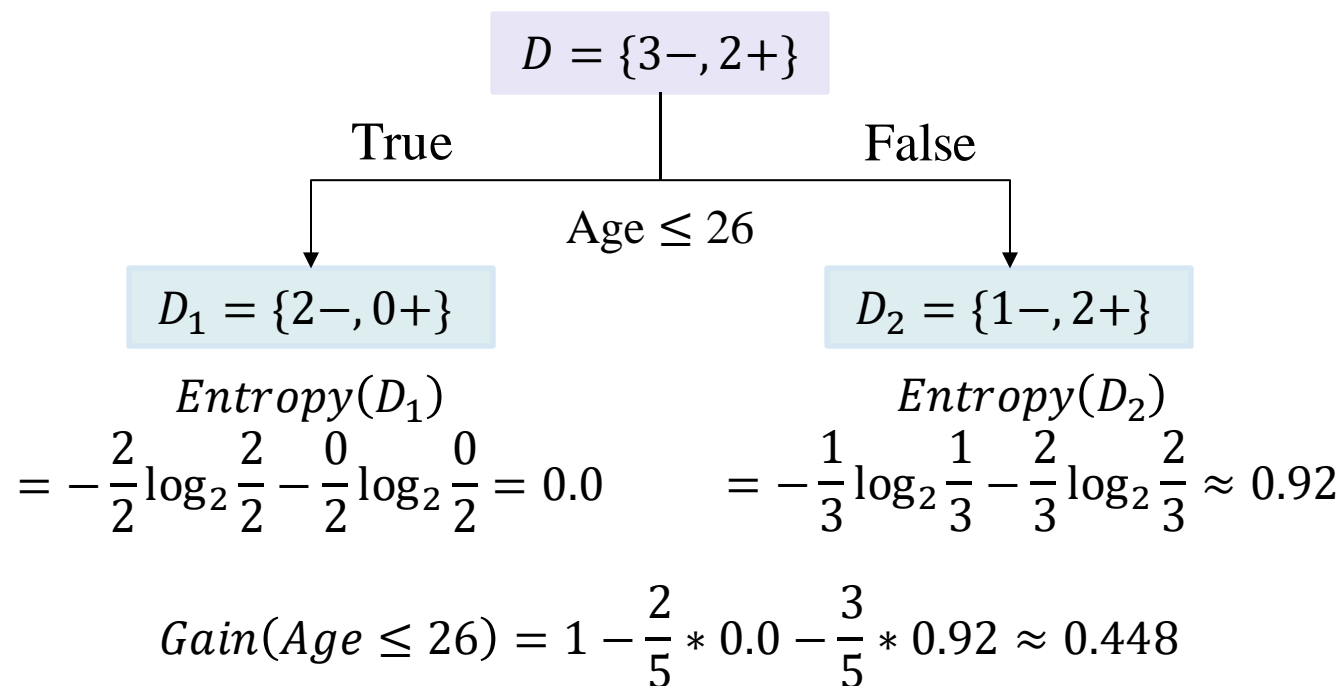
$$Entropy(D_i) = - \sum_{j=1}^c p_j \log_2 p_j$$

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Entropy / Gain



Age	Likes English	Likes AI	Raise Salary
23	0	0	0
25	1	1	0
27	1	0	1
29	0	1	1
29	0	0	0

# DT for Classification



## Constructing Decision Tree

$$Gain(D) = 1 - Entropy(D)$$

$$Entropy(D) = \frac{n_1}{n} Entropy(D_1) + \frac{n_2}{n_c} Entropy(D_2)$$

$$Entropy(D_i) = - \sum_{j=1} p_j \log_2 p_j$$

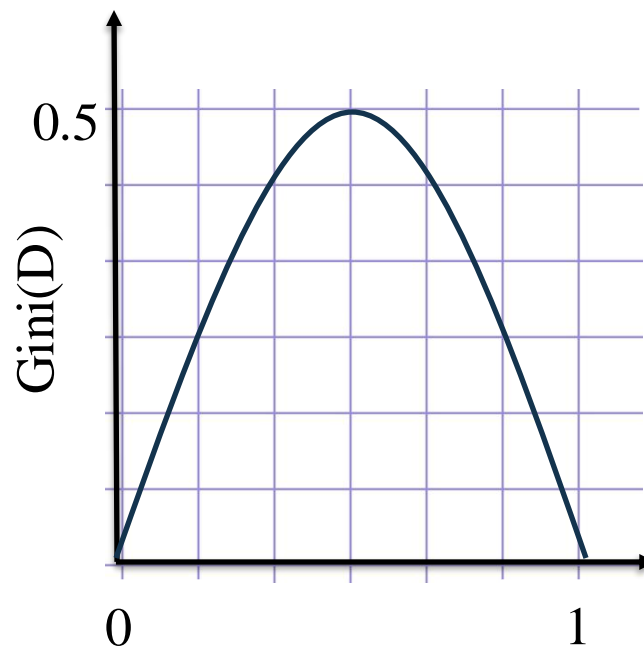
$$Gini(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

$$Gini(D_i) = 1 - \sum_{j=1}^c p_j^2$$

$$p_1 = p_2$$

$$\Rightarrow Gini(D) = 0.5$$

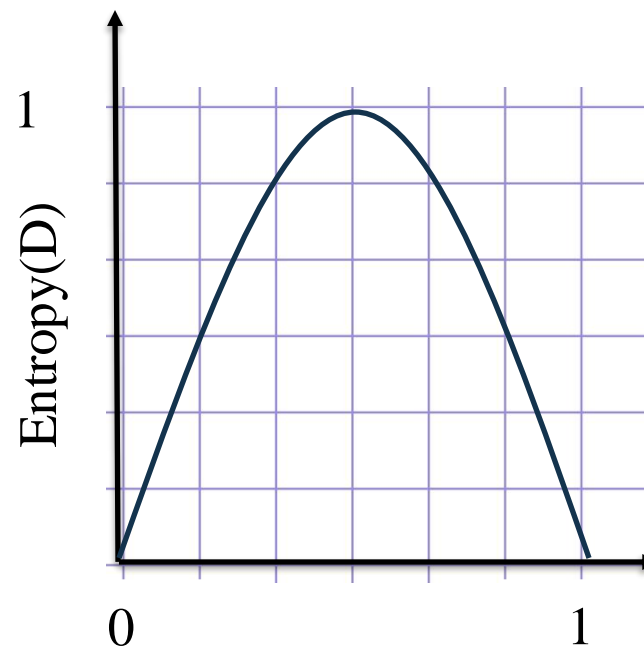
$$\Rightarrow Entropy(D) = 1.0$$



$$p_1 \text{ or } p_2 = 0$$

$$\Rightarrow Gini(D) = 0.0$$

$$\Rightarrow Entropy(D) = 0.0$$



# DT for Classification



## Decision Tree for IRIS Dataset

```
1 # Load the diabetes dataset
2 iris_X, iris_y = datasets.load_iris(return_X_y=True)
3
4 # Split train:test = 8:2
5 X_train, X_test, y_train, y_test = train_test_split(
6     iris_X, iris_y, test_size=0.2, random_state=42
7 )
8
9 # Train
10 dt_classifier = DecisionTreeClassifier()
11 dt_classifier.fit(X_train, y_train)
12
13 # Predict
14 y_pred = dt_classifier.predict(X_test)
15
16 # Evaluation
17 accuracy_score(y_test, y_pred)
```



QUIZ TIME

## SECTION 1

### Decision Tree

## SECTION 2

### DT for Classification

## SECTION 3

### DT for Regression

## SECTION 4

### Improved DT

Dataset D (n sample)



Dataset D<sub>1</sub>  
(n<sub>1</sub> sample)

Dataset D<sub>2</sub>  
(n<sub>2</sub> sample)

Age	Likes English	Likes AI	Salary
23	0	0	200
25	1	1	400
27	1	0	300
29	0	1	500
29	0	0	400

# DT for Regression



## Decision Tree



### Example Dataset

Age	Likes English	Likes AI	Salary
23	0	0	200
25	1	1	400
27	1	0	300
29	0	1	500
29	0	0	400

# DT for Regression



## Sum of Squared Errors (SSE)

$$SSE(D) = SSE(D_1) + SSE(D_2)$$

$$SSE(D_i) = \frac{1}{n_i} \sum_{j=1}^{n_i} (x_j - \bar{x}_i)^2$$

$x_i, \dots, x_{n_i}$  the items of  $D_i$

1. For each possible split, create child nodes and calculate the variance of each child node.
2. Calculate the variance of the split as the weighted average variance of child nodes.
3. Select the split with the lowest variance.

Dataset D (n sample)



Dataset D<sub>1</sub>  
(n<sub>1</sub> sample)

Dataset D<sub>2</sub>  
(n<sub>2</sub> sample)

Age	Likes English	Likes AI	Salary
23	0	0	<b>200</b>
25	1	1	<b>400</b>
27	1	0	<b>300</b>
29	0	1	<b>500</b>
29	0	0	<b>400</b>

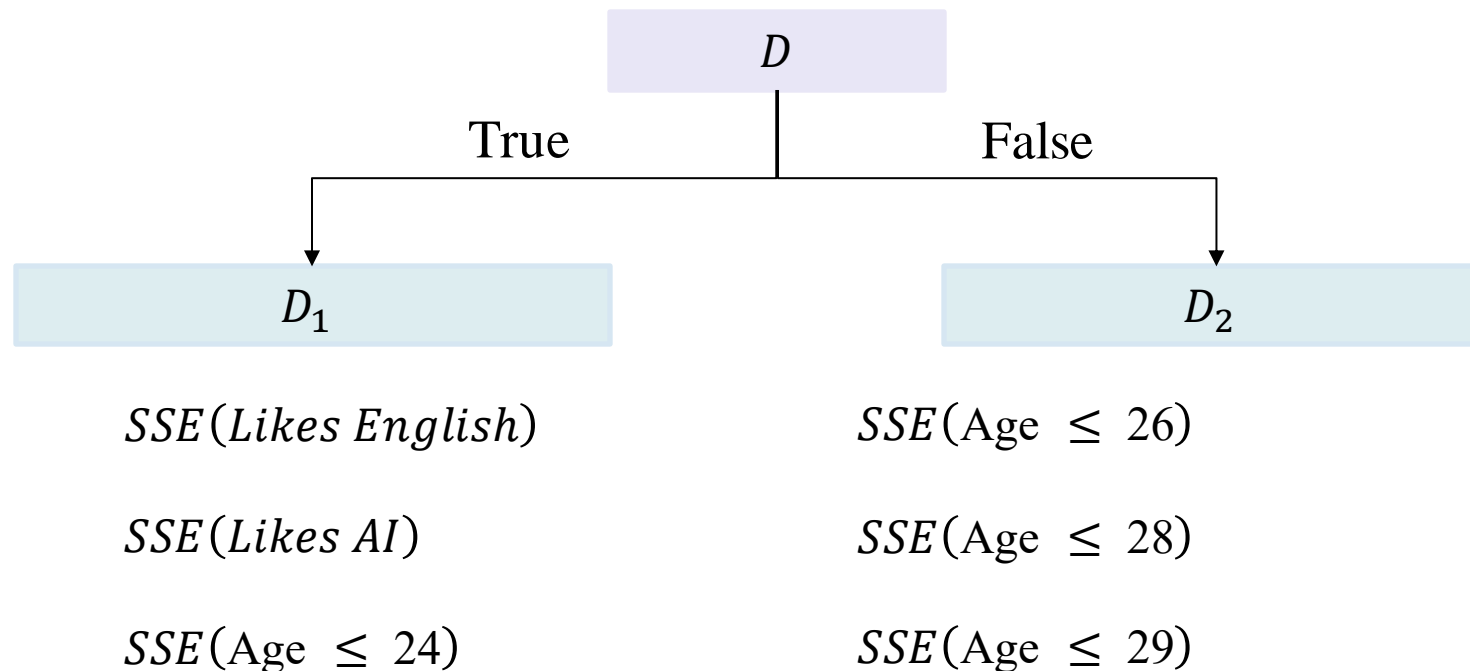
# DT for Regression



## Sum of Squared Errors (SSE)

$$SSE(D) = SSE(D_1) + SSE(D_2)$$

$$SSE(D_i) = \frac{1}{n_i} \sum_{j=1}^{n_i} (x_j - \bar{x}_i)^2$$



Age	Likes English	Likes AI	Salary
23	0	0	<b>200</b>
25	1	1	<b>400</b>
27	1	0	<b>300</b>
29	0	1	<b>500</b>
29	0	0	<b>400</b>

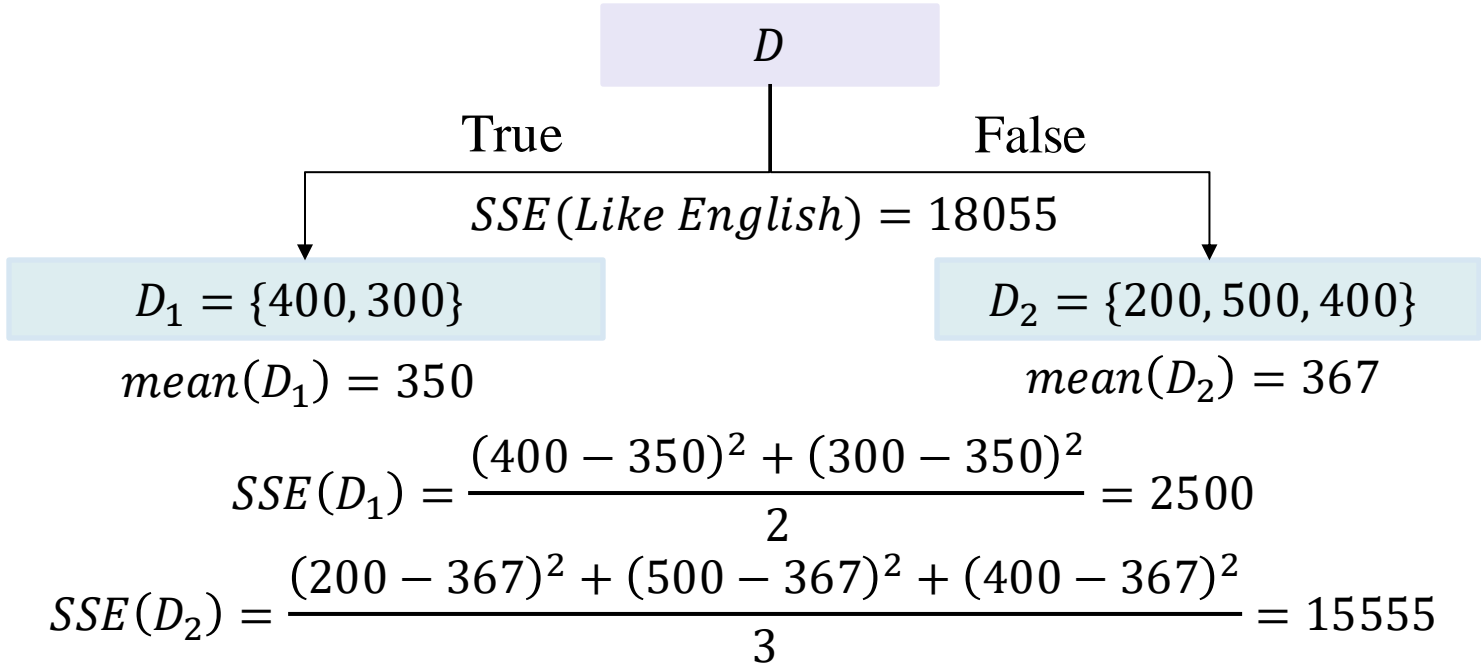
# DT for Regression



## Sum of Squared Errors (SSE)

$$SSE(D) = SSE(D_1) + SSE(D_2)$$

$$SSE(D_i) = \frac{1}{n_i} \sum_{j=1}^{n_i} (x_j - \bar{x}_i)^2$$



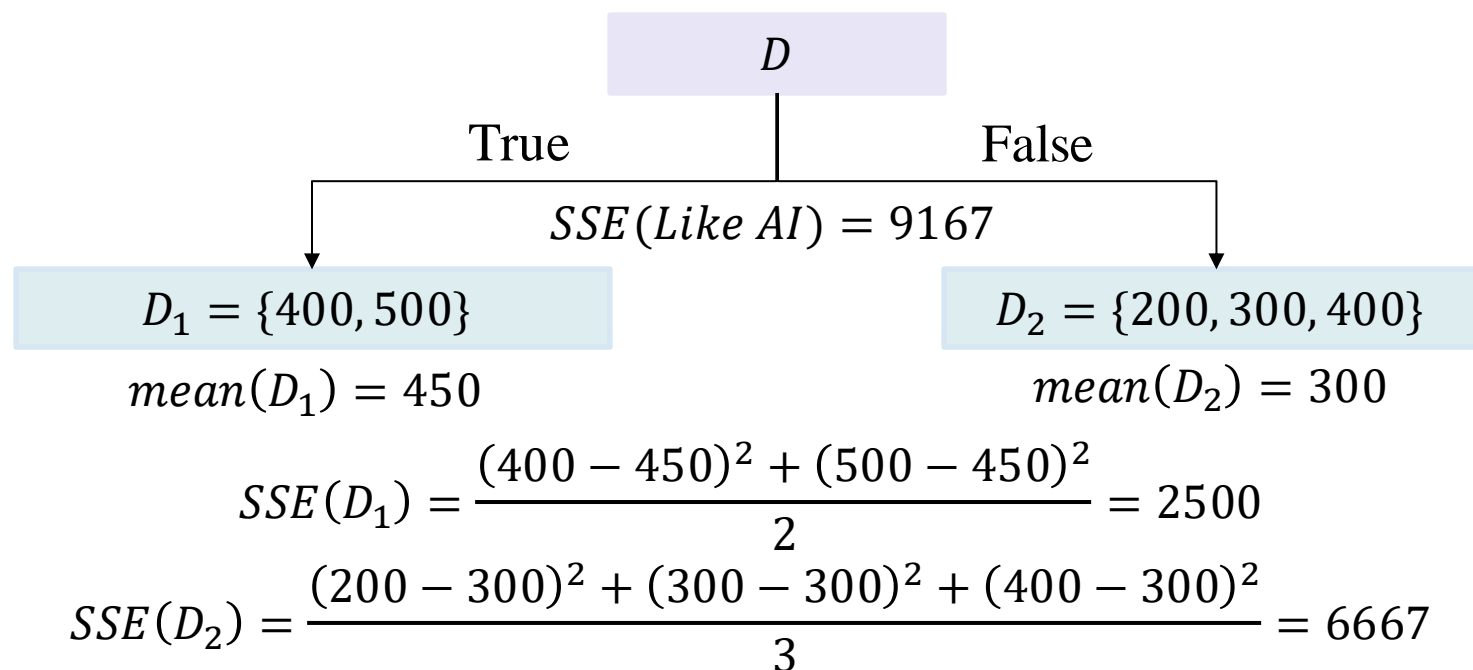
Age	Likes English	Likes AI	Salary
23	0	0	200
25	1	1	400
27	1	0	300
29	0	1	500
29	0	0	400

# DT for Regression

## ! Sum of Squared Errors (SSE)

$$SSE(D) = SSE(D_1) + SSE(D_2)$$

$$SSE(D_i) = \frac{1}{n_i} \sum_{j=1}^{n_i} (x_j - \bar{x}_i)^2$$



Age	Likes English	Likes AI	Salary
23	0	0	200
25	1	1	400
27	1	0	300
29	0	1	500
29	0	0	400

# DT for Regression

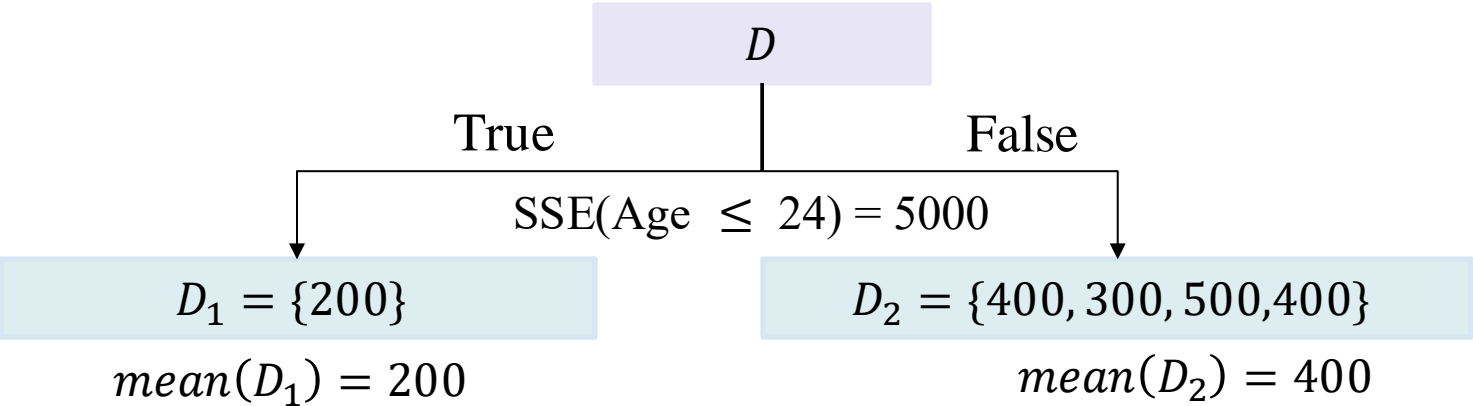


## Sum of Squared Errors (SSE)

$$SSE(D) = SSE(D_1) + SSE(D_2)$$

$$SSE(D_i) = \frac{1}{n_i} \sum_{j=1}^{n_i} (x_j - \bar{x}_i)^2$$

Numerical Feature	Ascending Ordering	Calculate mean	Determine Gini
-------------------	--------------------	----------------	----------------



$$SSE(D_1) = \frac{(200 - 200)^2}{1} = 0$$

$$SSE(D_2) = \frac{(400 - 400)^2 + (300 - 400)^2 + (500 - 400)^2 + (400 - 400)^2}{4} = 5000$$

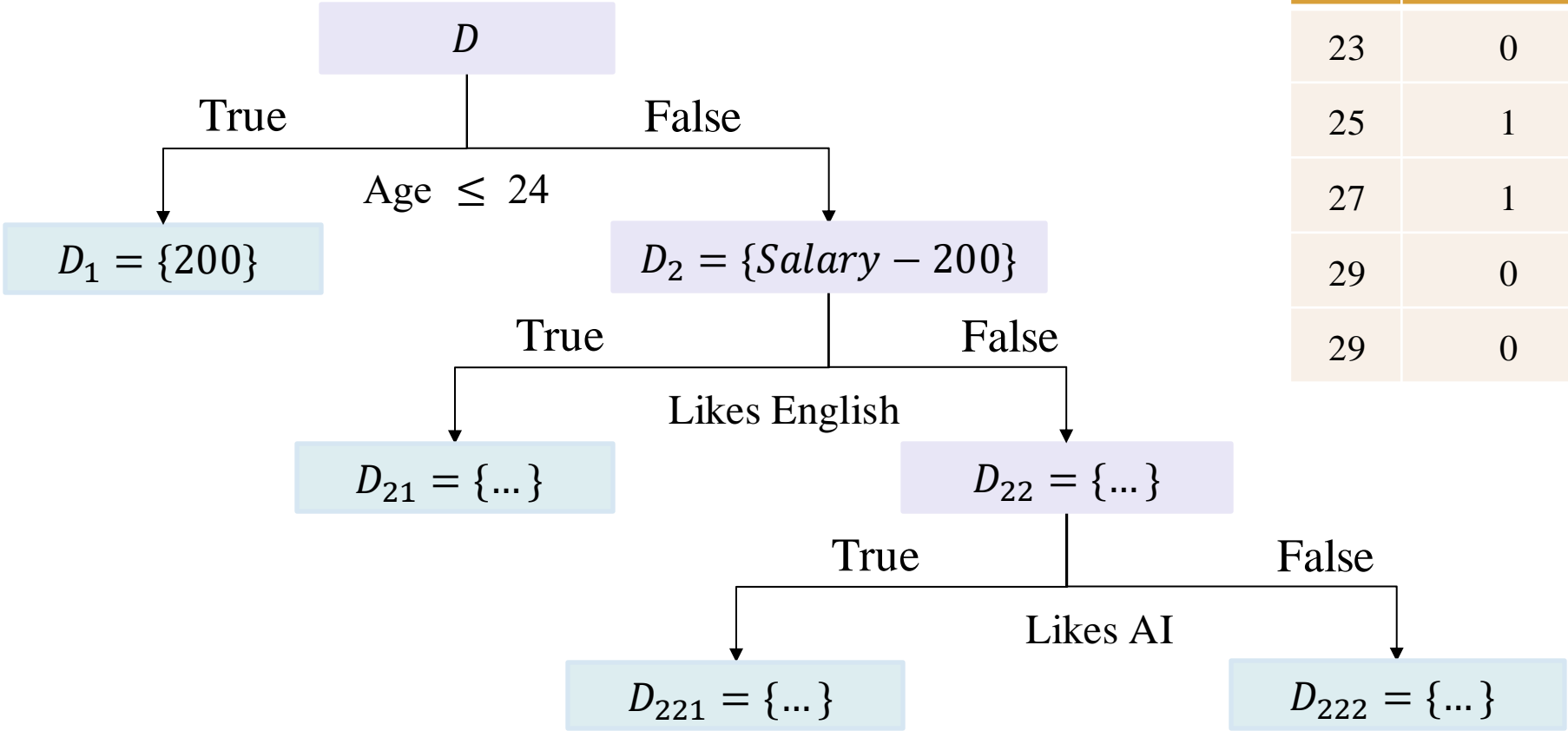
Age	Likes English	Likes AI	Salary
23	0	0	200
25	1	1	400
27	1	0	300
29	0	1	500
29	0	0	400



# DT for Regression



## Sum of Squared Errors (SSE)



Age	Likes English	Likes AI	Salary
23	0	0	200
25	1	1	400
27	1	0	300
29	0	1	500
29	0	0	400

# DT for Regression



## Decision Tree for CPU Machine Regression Dataset

```
1 # Load dataset
2 from sklearn.datasets import fetch_openml
3
4 machine_cpu = fetch_openml(name='machine_cpu')
5
6 machine_data = machine_cpu.data
7 machine_labels = machine_cpu.target
8
9 X_train, X_test, y_train, y_test = train_test_split(
10     machine_data, machine_labels, test_size=0.2, random_state=20
11 )
12
13 # train
14 tree_reg = DecisionTreeRegressor()
15 tree_reg.fit(X_train, y_train)
16
17 # predict
18 y_pred = tree_reg.predict(X_test)
19
20 # evaluation
21 mean_squared_error(y_test, y_pred)
```

## SECTION 1

## Decision Tree

## SECTION 2

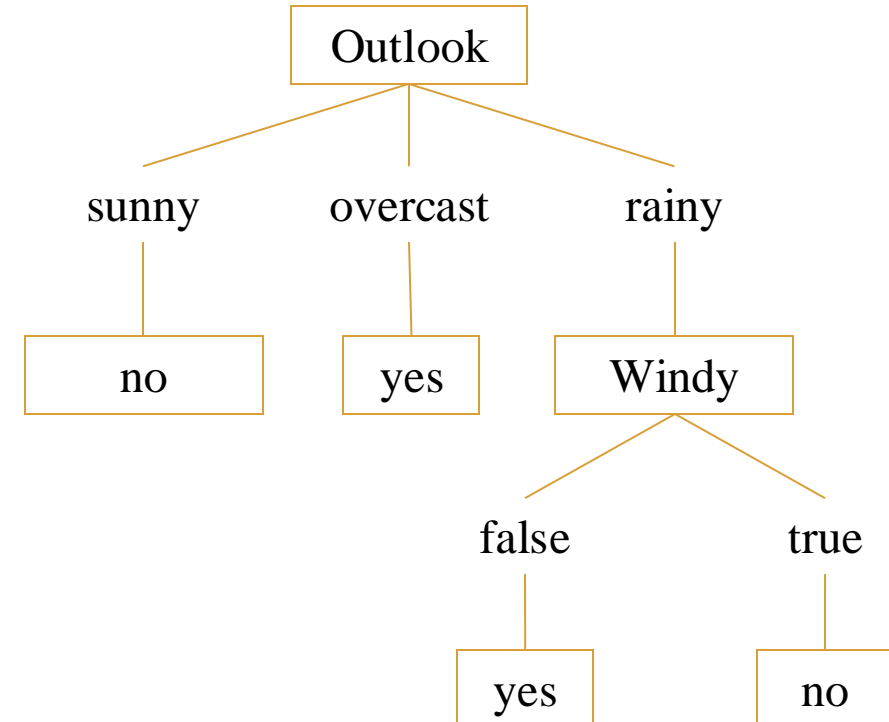
## DT for Classification

## SECTION 3

## DT for Regression

## SECTION 4

## Improved DT

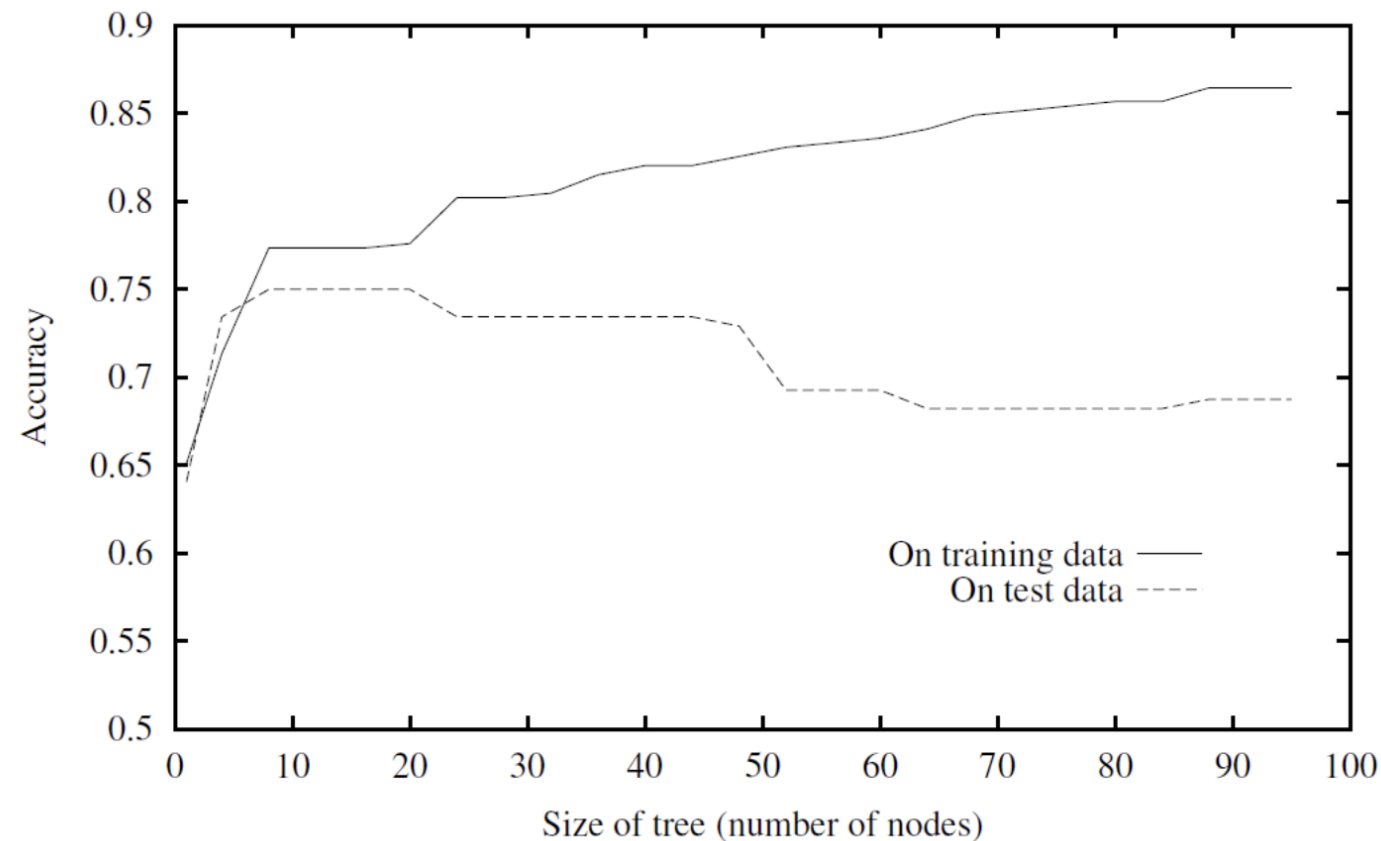


# Improved Decision Tree



## Over-fitting

- Desired: a DT that is not too big in size, yet fits the training data reasonably



[Source](#)

# Improved Decision Tree



## Over-fitting

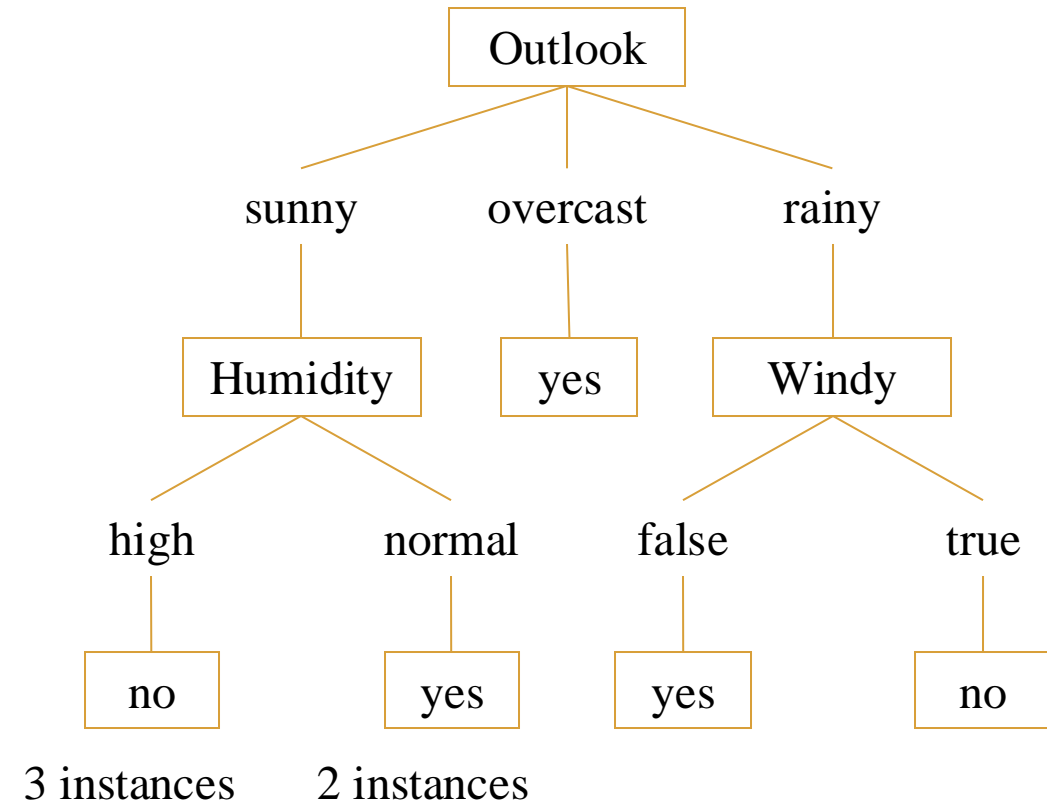
- Mainly two approaches
  - ❑ Prune while building the tree (Stopping Early)
  - ❑ Prune after building the tree (Post-Pruning)
- Criteria for judging which nodes could potentially be pruned: evaluate validation set
  - ❑ Reduced-error pruning
  - ❑ Rule post-pruning

# Improved Decision Tree



## Reduced-Error Pruning

- Pruning a decision node  $d$  consists of:
  - ❑ Removing the subtree rooted at  $d$
  - ❑ Making  $d$  a leaf node
  - ❑ Assigning  $d$  the most common classification of the training instances associated with  $d$ .



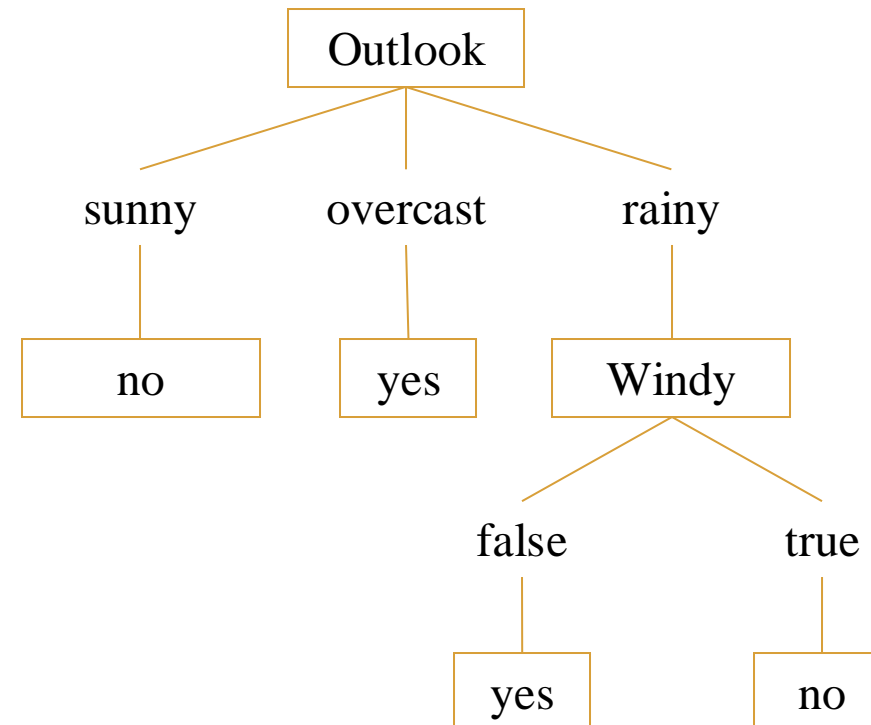
Accuracy of the tree on the validation set is 90%.

# Improved Decision Tree



## Reduced-Error Pruning

- Pruning a decision node  $d$  consists of:
  - ❑ Removing the subtree rooted at  $d$
  - ❑ Making  $d$  a leaf node
  - ❑ Assigning  $d$  the most common classification of the training instances associated with  $d$ .
  
- Do until further pruning is harmful:
  - ❑ Evaluate impact on validation set of pruning each possible node
  - ❑ Greedily remove the one that most improves validation set accuracy.



Accuracy of the tree on the validation set is 92.4%.

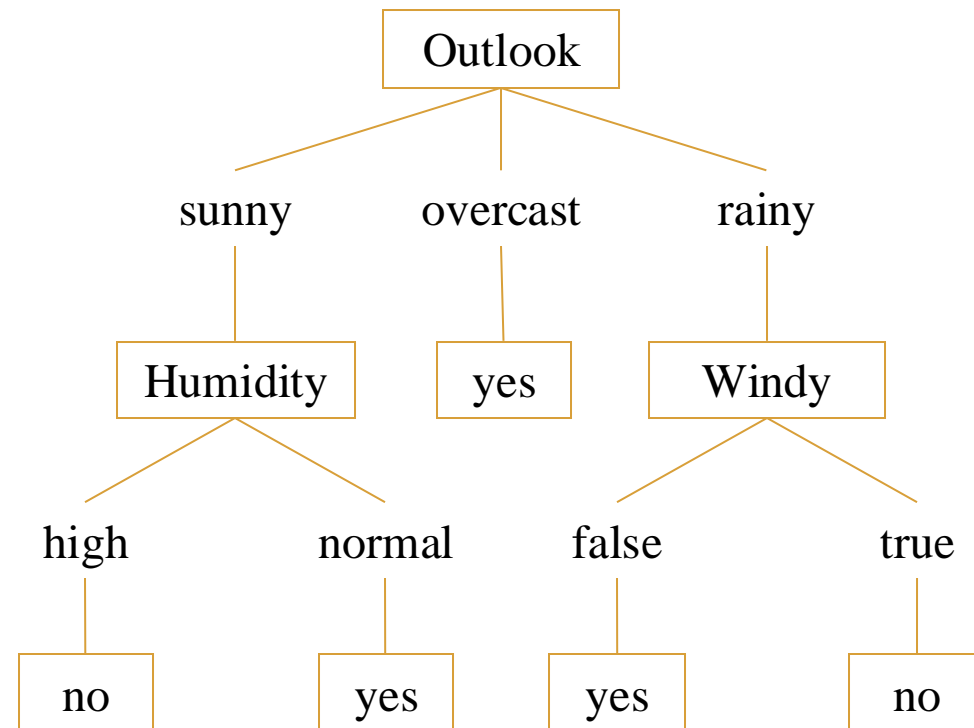
# Improved Decision Tree



## Rule Post-Pruning

- Convert tree to equivalent set of rules.
- Prune each rule independently of others
- Sort final rules by their estimated accuracy and consider them in this sequence when classifying subsequent instances.

IF (Outlook = Sunny) & (Humidity = High)  
THEN PlayTennis = No  
IF (Outlook = Sunny) & (Humidity = Normal)  
THEN PlayTennis = Yes  
.....



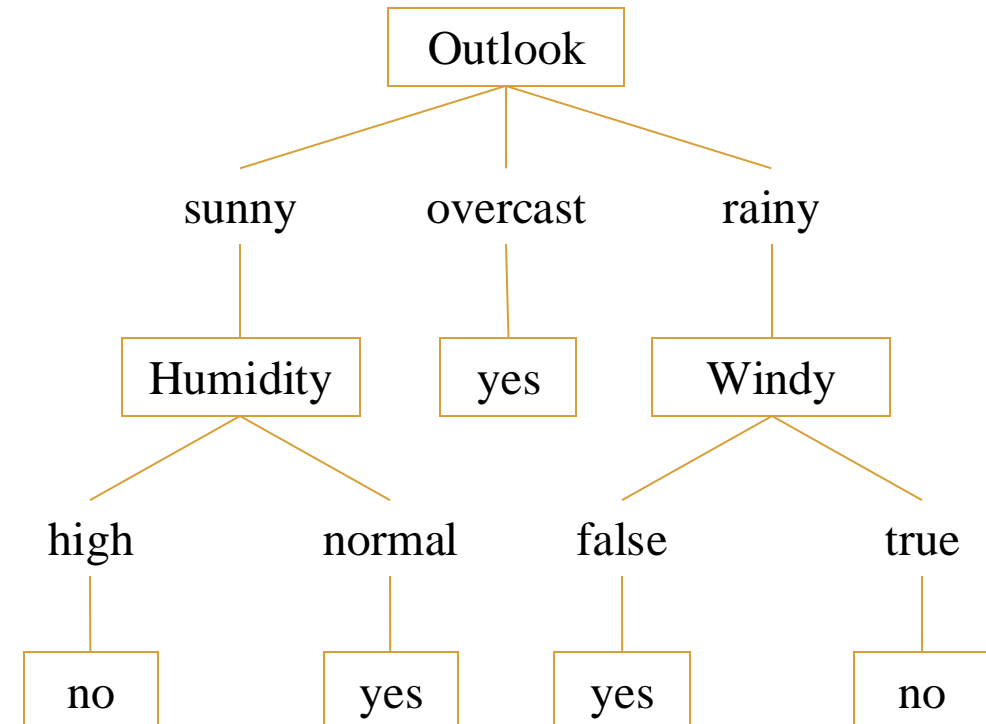


# Improved Decision Tree



## Missing Attribute Values

- Assign most common value of A among other instances belonging to the same concept.
- If node n test the attribute A, assign most common value of A among other instances sorted to node n.
- If node n tests the attribute A, assign a probability to each of possible values of A. These probabilities are estimated based on the observed frequencies of the values of A (Based on information gain).



Test = {Outlook=Sunny, **Humidity=?**, Windy=True}

# Summary

## Decision Tree

- ❖ Introduction
- ❖ Regression & Classification Problem
- ❖ Terminology

## Decision Tree for Regression

- ❖ Variance
- ❖ Sum of Squared Errors (SSE)
- ❖ Based on Sklearn Library

## Decision Tree for Classification

- ❖ Constructing Decision Tree: ID3
- ❖ Gini Impurity
- ❖ Entropy
- ❖ Information Gain
- ❖ Based on Sklearn library

## Improved Decision Tree

- ❖ Overfitting
- ❖ Stopping Early
- ❖ Post-Pruning (Reduced-Error & Rule Post-Pruning)
- ❖ Missing Attribute Values



AI VIET NAM

@aivietnam.edu.vn

# Thanks!

## Any questions?