

## Module 03 – Extra Class

# RANDOM FOREST

Nguyen Quoc Thai

# Objectives

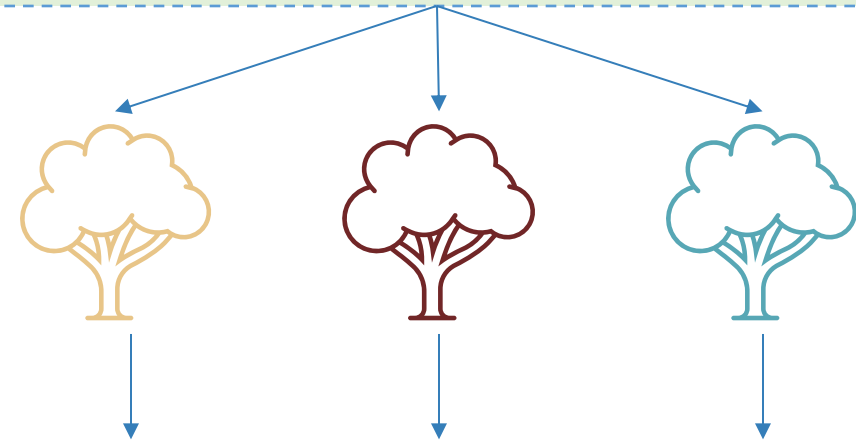
## Decision Tree Review

- ❖ Decision Tree
- ❖ Decision Tree for Classification
- ❖ Decision Tree for Regression
- ❖ IRIS Classification
- ❖ Salary Prediction



## Random Forest

- ❖ Decision Tree
- ❖ Random Forest
  - Bootstrap Sample
  - Majority Voting / Averaging
- ❖ IRIS Classification
- ❖ Salary Prediction



# Outline

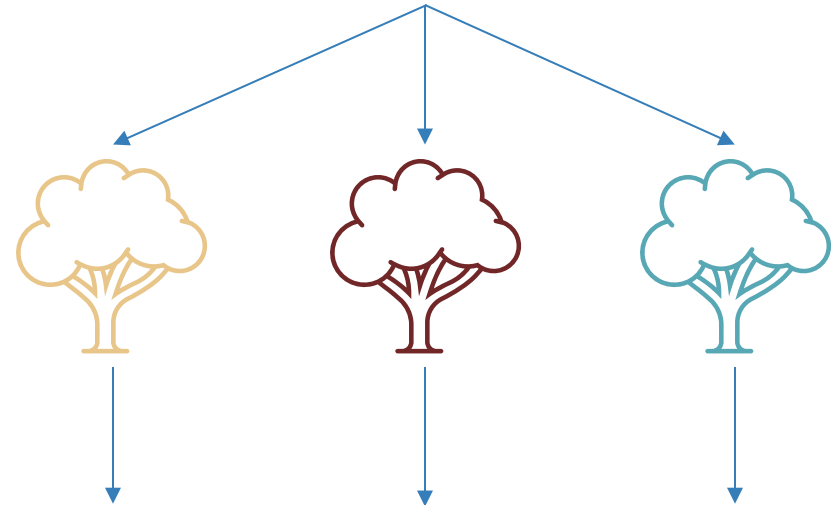
## SECTION 1

### Decision Tree Review



## SECTION 2

### Random Forest

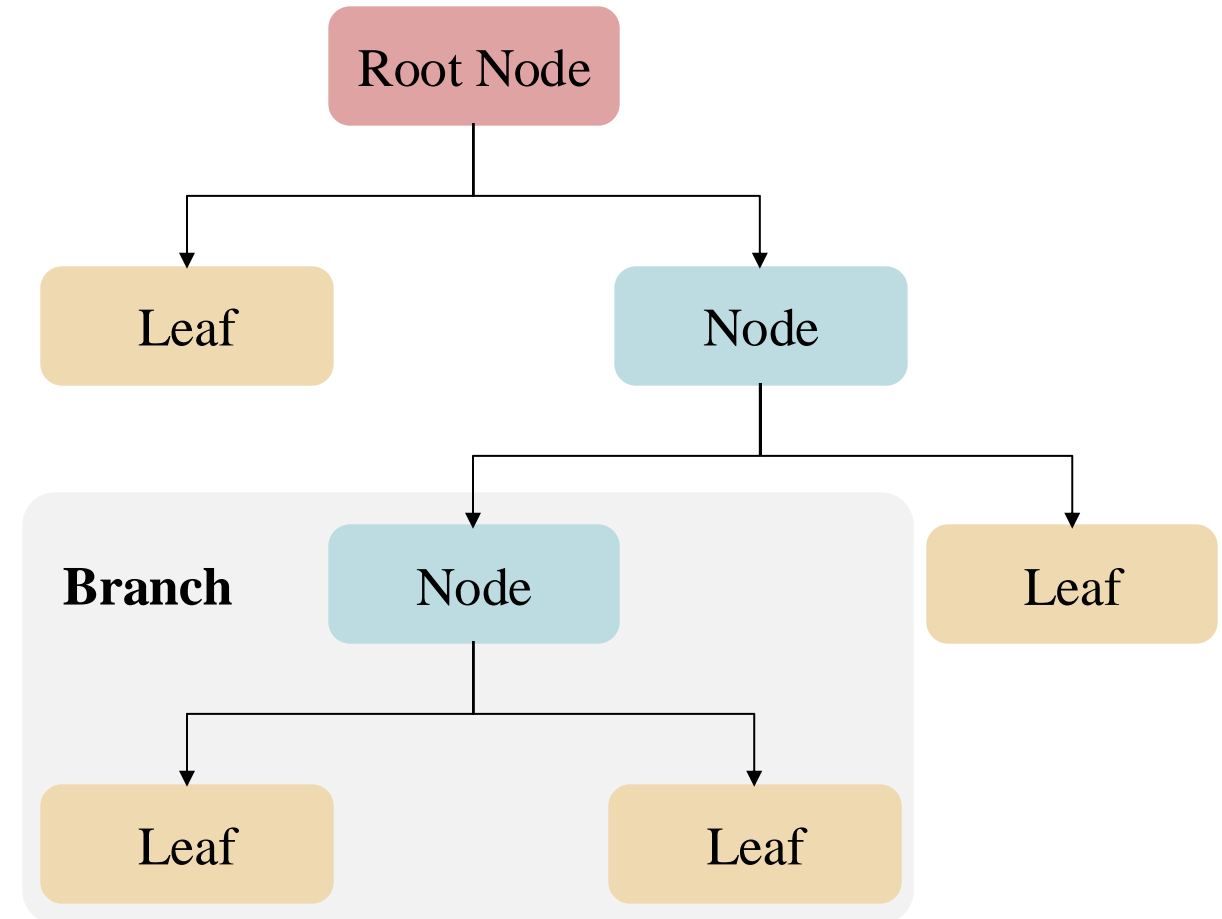


# Decision Tree Review



## Decision Tree

- ❖ **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 Review



## Decision Tree for Classification

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

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

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

$$Gini(D) = 1 - \left(\frac{3}{3+3}\right)^2 - \left(\frac{3}{3+3}\right)^2 = \frac{1}{2}$$

Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

# Decision Tree Review



## Decision Tree for Classification

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

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

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini

Petal_Length	
0.9	0.95
1	
1.2	
1.3	
1.7	
1.8	1.75

$Gini(Length \leq 0.95)$

$Gini(Length \leq 1.1)$

$Gini(Length \leq 1.25)$

$Gini(Length \leq 1.5)$

$Gini(Length \leq 1.75)$

Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

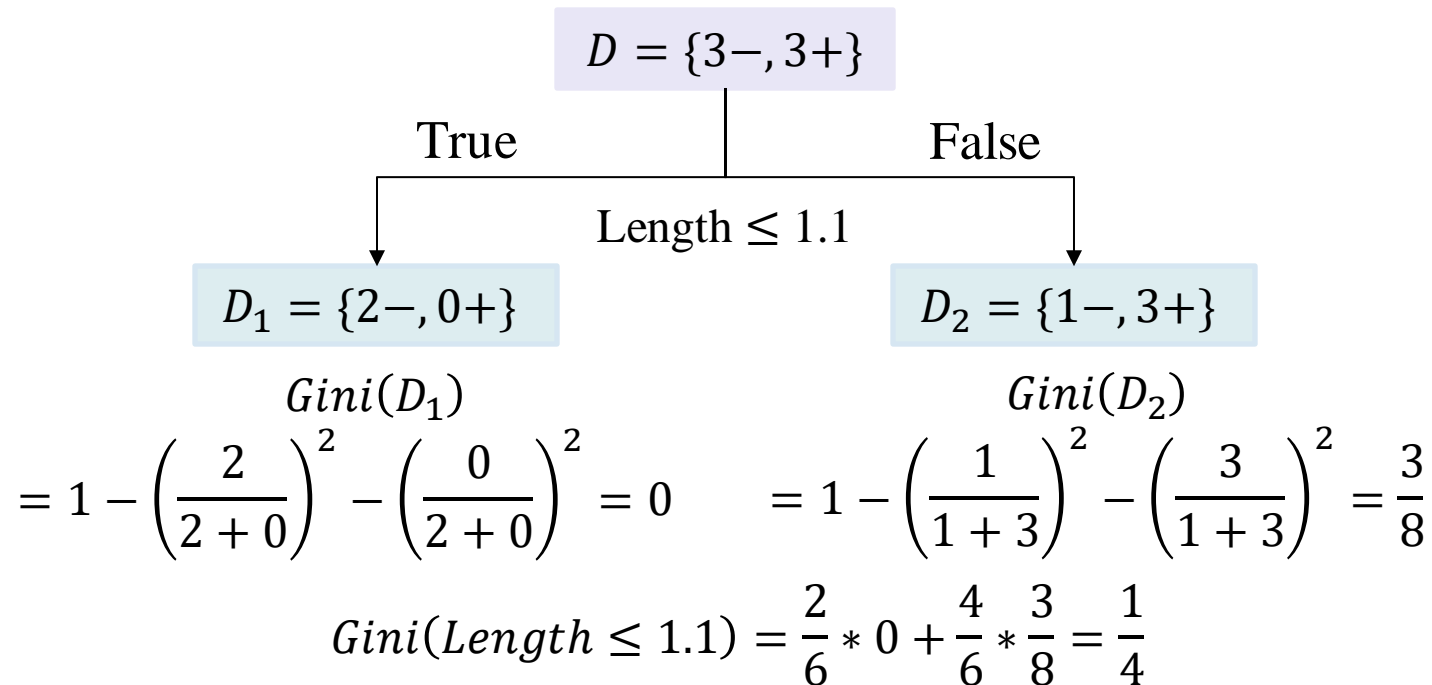
# Decision Tree Review



## Decision Tree for Classification

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

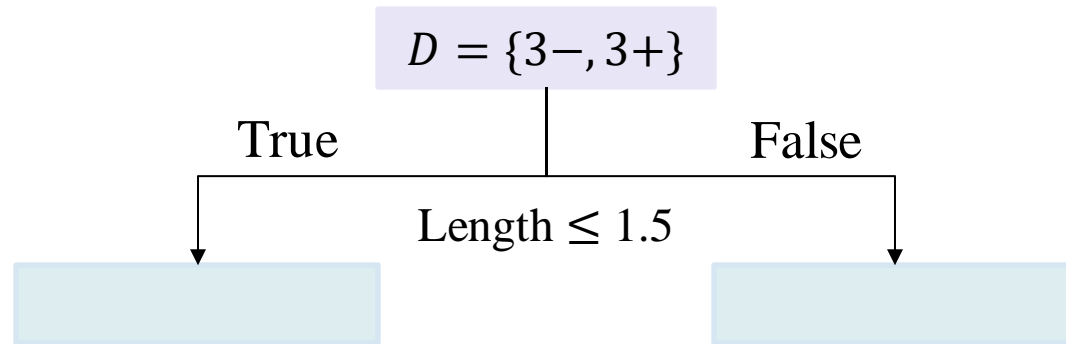
# Decision Tree Review



## Decision Tree for Classification

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1



# Decision Tree Review



## Decision Tree for Classification

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

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

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini

Petal_Width		
0.2	0.35 0.55 0.65 0.8 1.1	$Gini(Width \leq 0.35)$
0.5		$Gini(Width \leq 0.55)$
0.6		$Gini(Width \leq 0.65)$
0.7		$Gini(Width \leq 0.8)$
0.9		$Gini(Width \leq 1.1)$
1.3		

Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

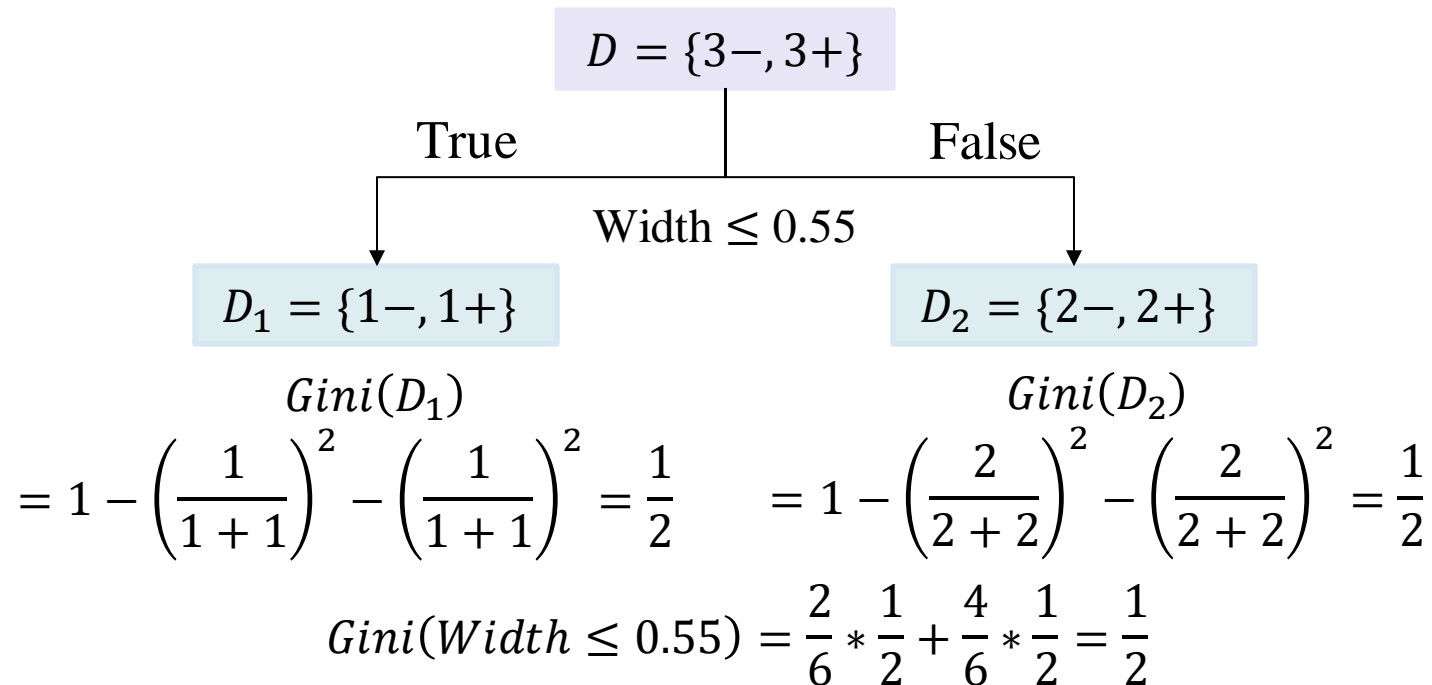
# Decision Tree Review



## Decision Tree for Classification

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

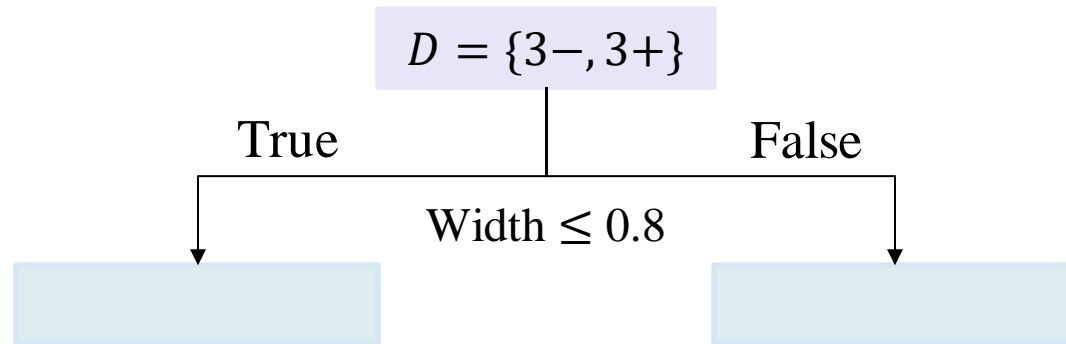
# Decision Tree Review



## Decision Tree for Classification

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

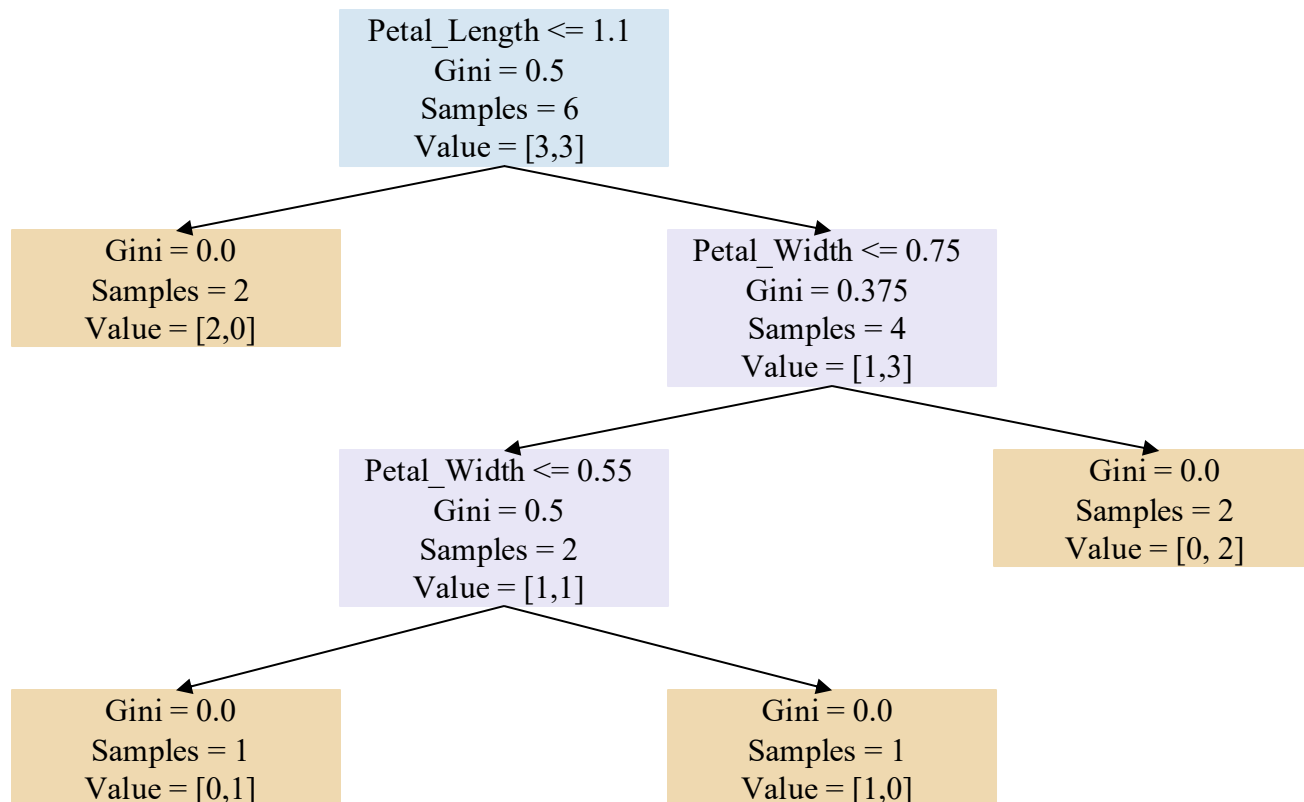
# Decision Tree Review



## Decision Tree for Classification

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

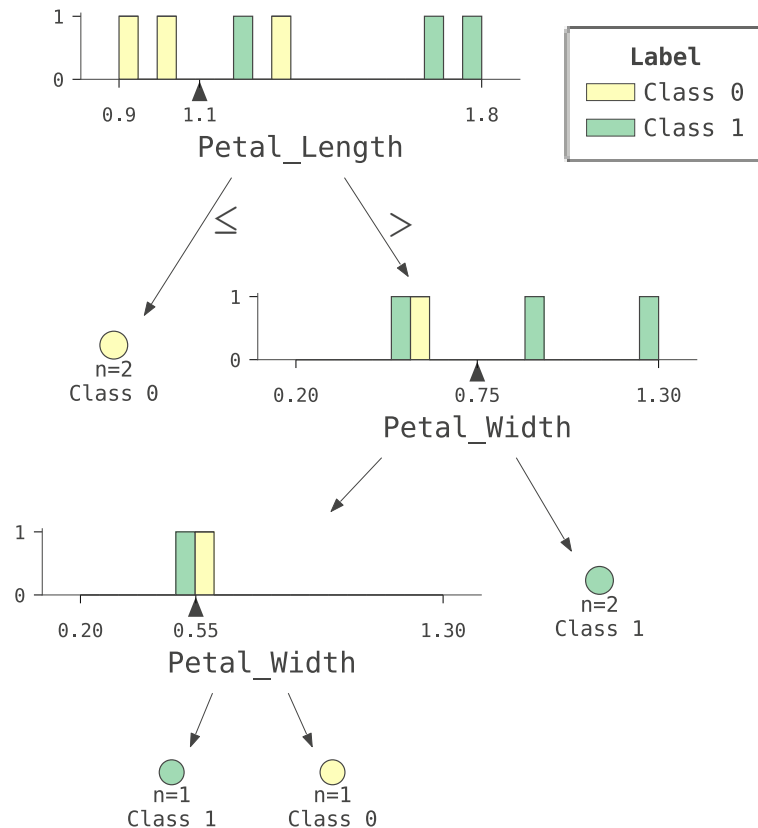
# Decision Tree Review



## Decision Tree for Classification

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

# Decision Tree Review



## Decision Tree for Classification

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

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

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

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

$$Entropy(D) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1$$

Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

# Decision Tree Review



## Decision Tree for Classification

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

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

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

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini

Petal_Length		
0.9	0.95 1.1 1.25 1.5 1.75	Entropy( $Length \leq 0.95$ ) $Gain(Length \leq 0.95)$ Entropy( $Length \leq 1.1$ ) $Gain(Length \leq 1.1)$ Entropy( $Length \leq 1.25$ ) $Gain(Length \leq 1.25)$ Entropy( $Length \leq 1.5$ ) $Gain(Length \leq 1.5)$ Entropy( $Length \leq 1.75$ ) $Gain(Length \leq 1.75)$
1		
1.2		
1.3		
1.7		
1.8		

Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

# Decision Tree Review

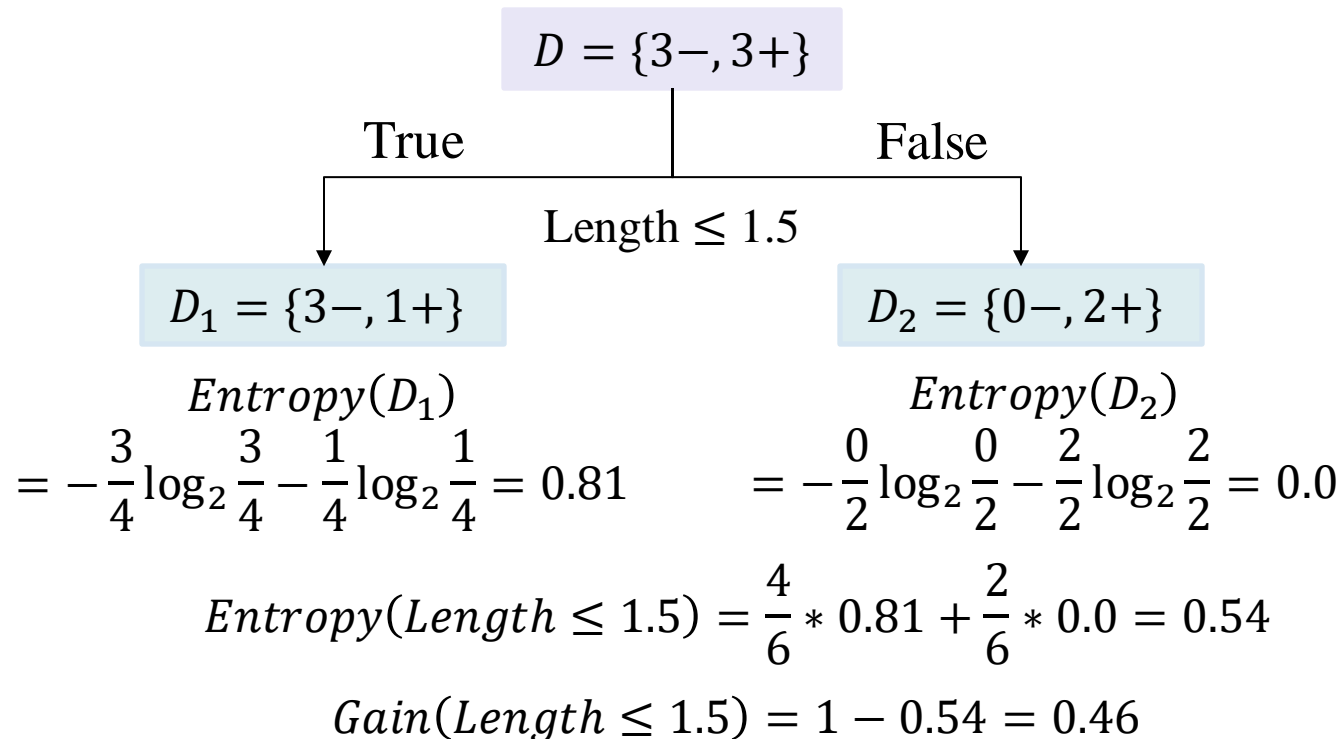


## Decision Tree for Classification

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

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1



# Decision Tree Review

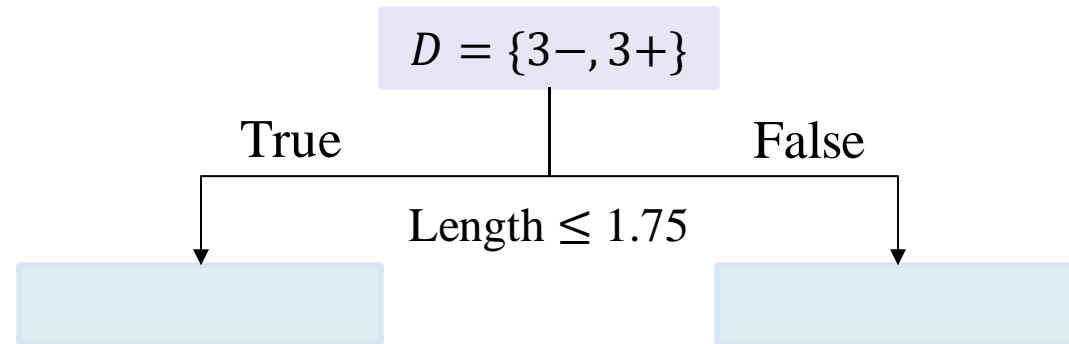


## Decision Tree for Classification

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

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

# Decision Tree Review



## Decision Tree for Classification

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

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

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

**Numerical Feature**

Ascending Ordering

Calculate mean

Determine Gini

Petal_Width				Petal_Length	Petal_Width	Label
0.2	0.35	Entropy( $Width \leq 0.35$ )	Gain( $Width \leq 0.35$ )	1	0.2	0
0.5	0.55	Entropy( $Width \leq 0.55$ )	Gain( $Width \leq 0.55$ )	1.3	0.6	0
0.6	0.65	Entropy( $Width \leq 0.65$ )	Gain( $Width \leq 0.65$ )	0.9	0.7	0
0.7	0.8	Entropy( $Width \leq 0.8$ )	Gain( $Width \leq 0.8$ )	1.7	0.5	1
0.9	1.1	Entropy( $Width \leq 1.1$ )	Gain( $Width \leq 1.1$ )	1.8	0.9	1
1.3				1.2	1.3	1

# Decision Tree Review

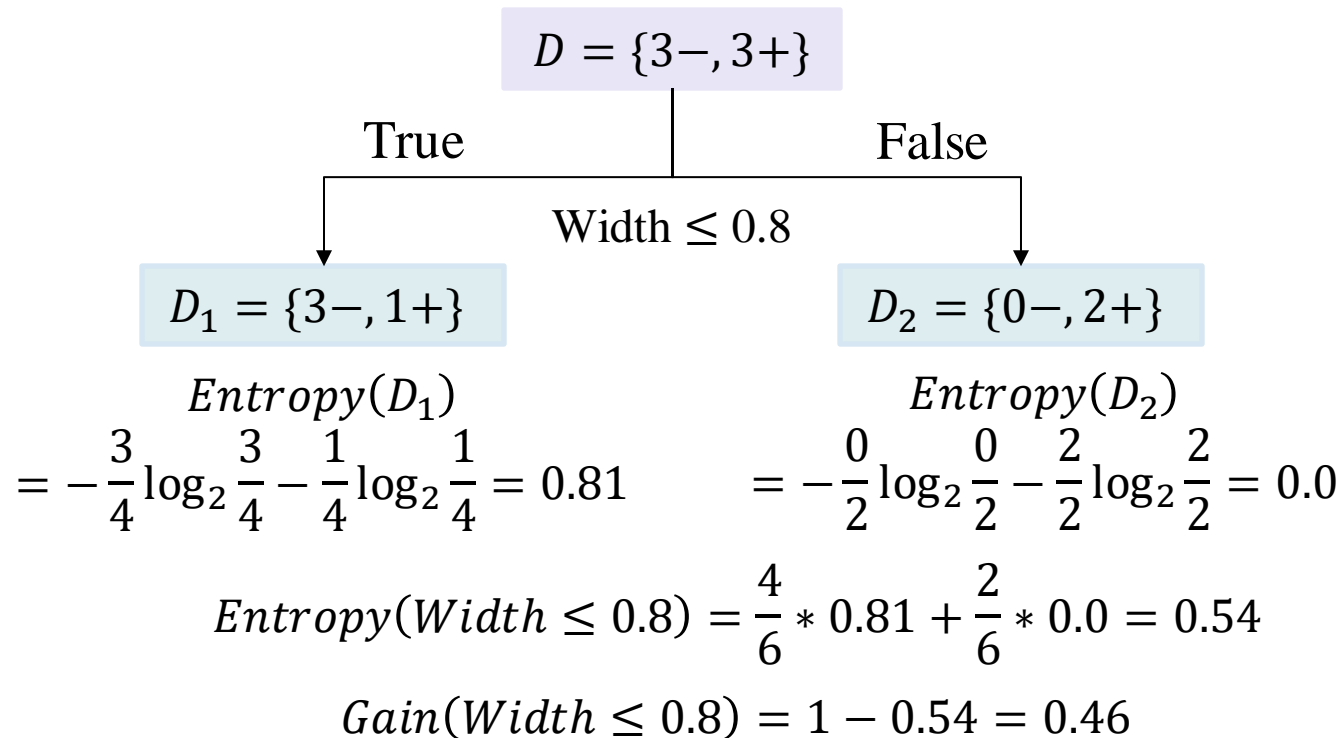


## Decision Tree for Classification

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

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

# Decision Tree Review

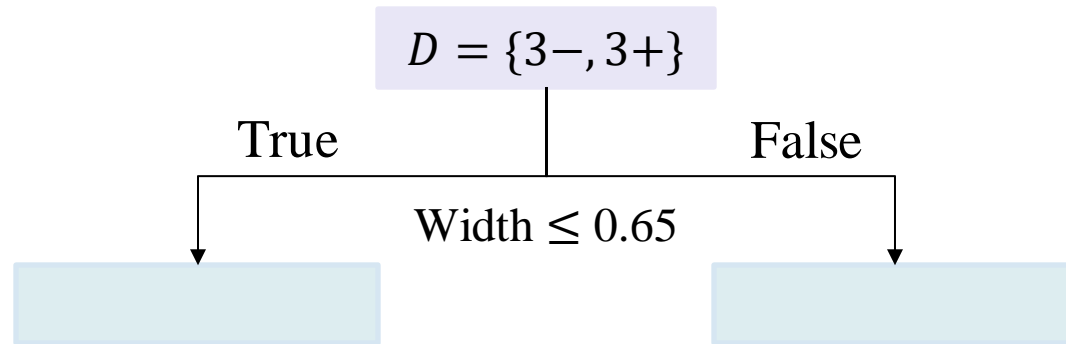


## Decision Tree for Classification

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

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

# Decision Tree Review

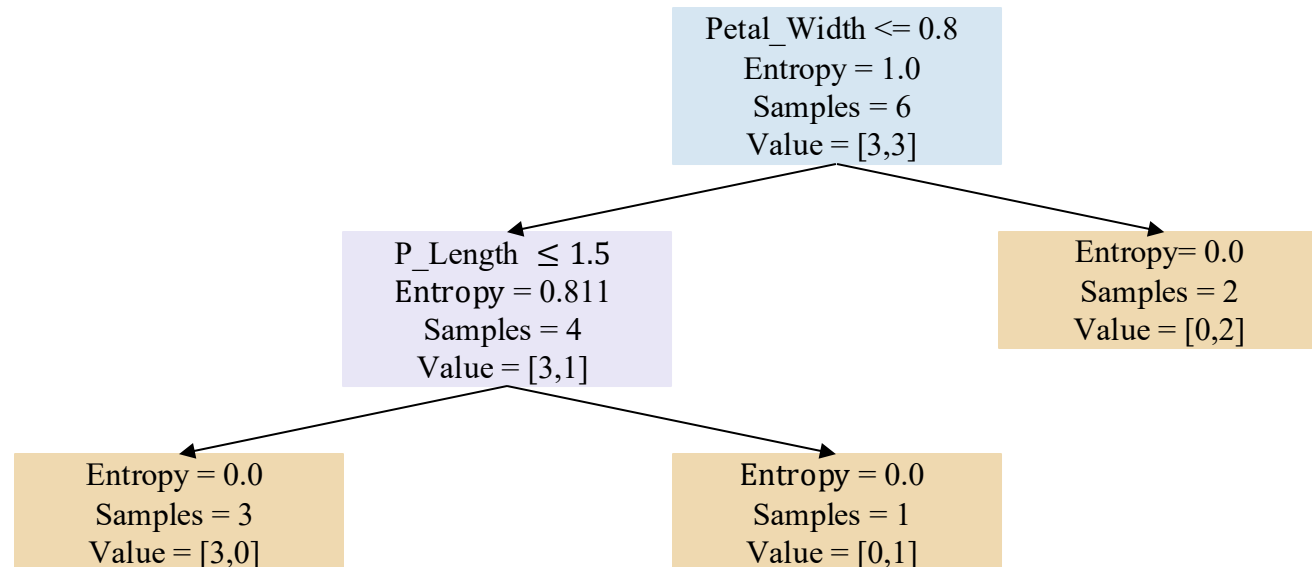


## Decision Tree for Classification

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

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

# Decision Tree Review

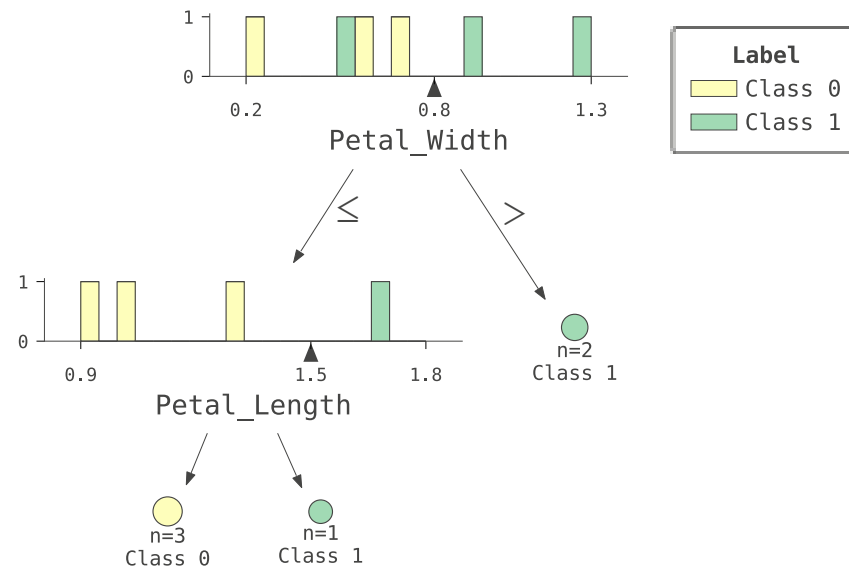


## Decision Tree for Classification

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

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

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



Petal_Length	Petal_Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

# Decision Tree Review



## Decision Tree for Regression

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

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

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

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

Experience	
1.5	2.0
2.5	
4.0	
5.5	
	3.25
	4.75

$$SSE(\text{Experience} \leq 2.0)$$

$$SSE(\text{Experience} \leq 3.25)$$

$$SSE(\text{Experience} \leq 4.75)$$

$$MSE(\text{Experience} \leq 2.0)$$

$$MSE(\text{Experience} \leq 3.25)$$

$$MSE(\text{Experience} \leq 4.75)$$

Experience	Salary
1.5	0
2.5	0
4.0	55
5.5	83

# Decision Tree Review



## Decision Tree for Regression

$$\begin{aligned} SSE(D) &= SSE(D_1) + SSE(D_2) \\ MSE(D) &= MSE(D_1) + MS(D_2) \end{aligned}$$

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

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

$$D = \{0, 0, 55, 83\}$$

$$\text{mean}(D) = 34.5$$

$$SSE(D_1) = (0 - 34.5)^2 + (0 - 34.5)^2 + (55 - 34.5)^2 + (83 - 34.5)^2 = 5153$$

$$MSE(D_1) = \frac{(0 - 34.5)^2 + (0 - 34.5)^2 + (55 - 34.5)^2 + (83 - 34.5)^2}{4} = 1288.25$$

Experience	Salary
1.5	0
2.5	0
4.0	55
5.5	83



# Decision Tree Review



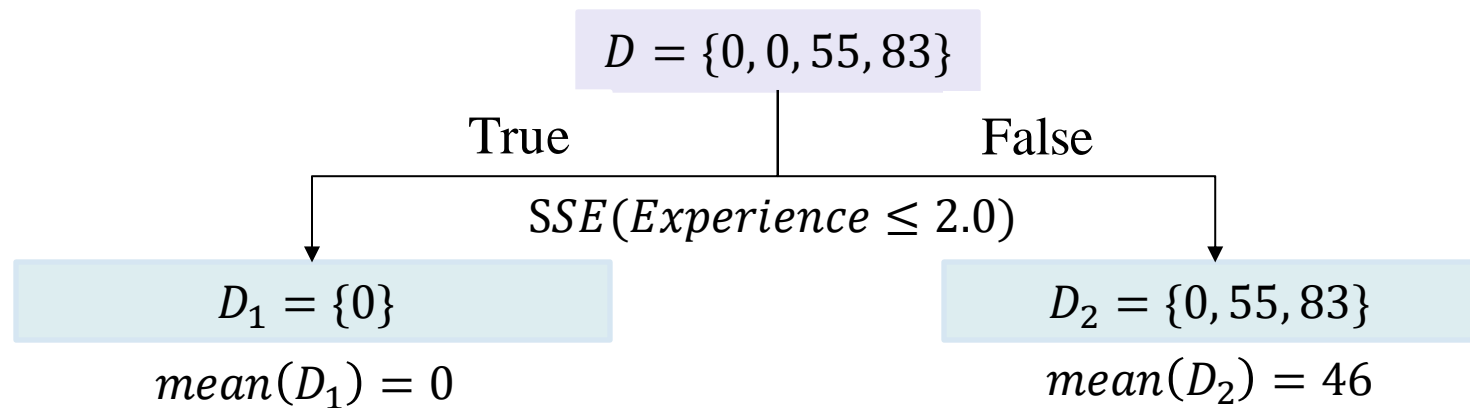
## Decision Tree for Regression

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

$$MSE(D) = MSE(D_1) + MS(D_2)$$

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

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



$$SSE(D_1) = (0 - 0)^2 = 0$$

$$SSE(D_2) = (0 - 46)^2 + (55 - 46)^2 + (83 - 46)^2 = 1450$$

$$SSE(Experience \leq 2.0) = 1450$$

Experience	Salary
1.5	0
2.5	0
4.0	55
5.5	83

# Decision Tree Review

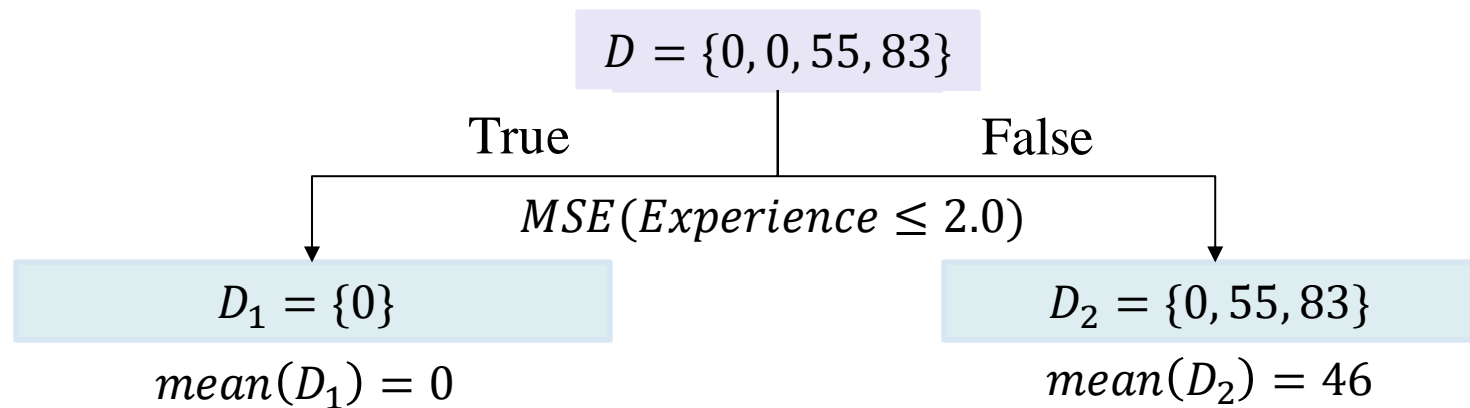
## ! Decision Tree for Regression

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

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

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

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



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

$$MSE(D_2) = \frac{(0 - 46)^2 + (55 - 46)^2 + (83 - 46)^2}{3} = 483$$

$$MSE(Experience \leq 2.0) = 483$$

Experience	Salary
1.5	0
2.5	0
4.0	55
5.5	83

# Decision Tree Review

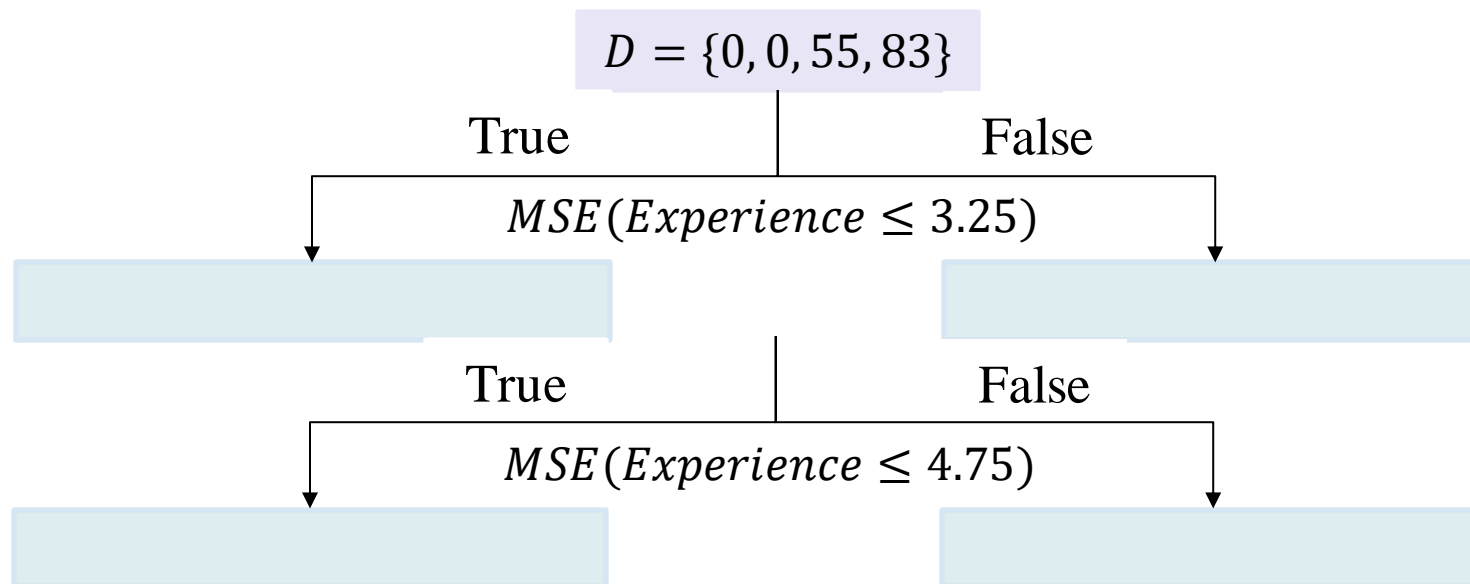
## ! Decision Tree for Regression

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

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

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

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



Experience	Salary
1.5	0
2.5	0
4.0	55
5.5	83

# Decision Tree Review



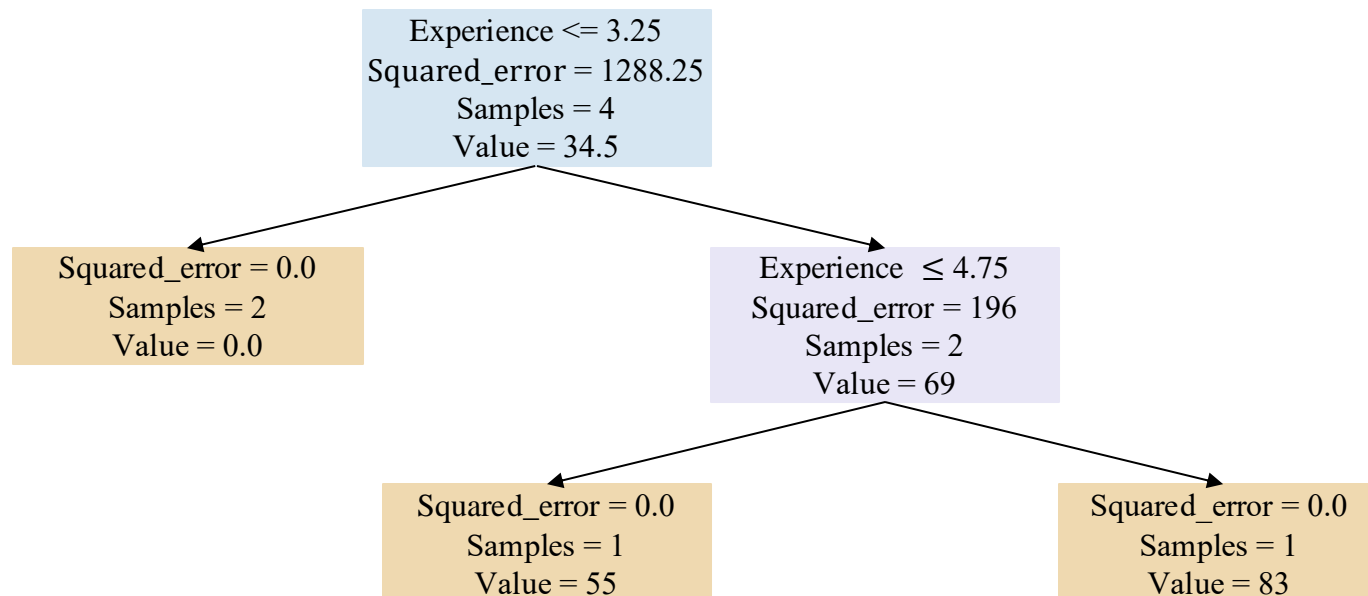
## Decision Tree for Regression

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

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

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

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



Experience	Salary
1.5	0
2.5	0
4.0	55
5.5	83

# Decision Tree Review



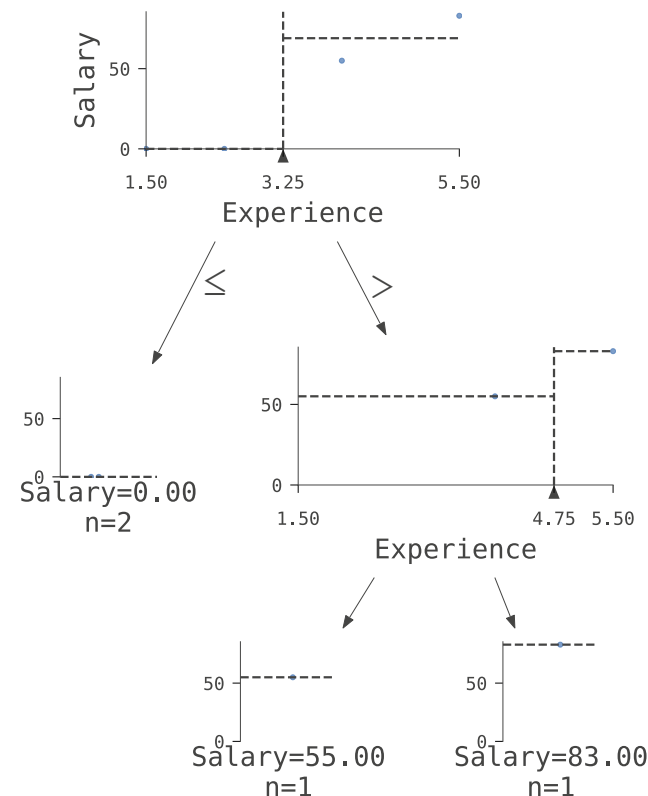
## Decision Tree for Regression

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

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

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

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



Experience	Salary
1.5	0
2.5	0
4.0	55
5.5	83

# Outline

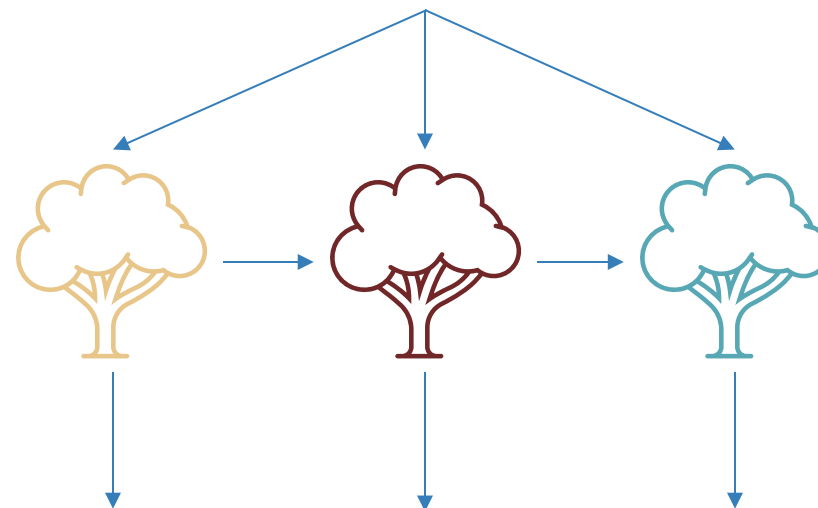
## SECTION 1

### Decision Tree Review



## SECTION 2

### Random Forest

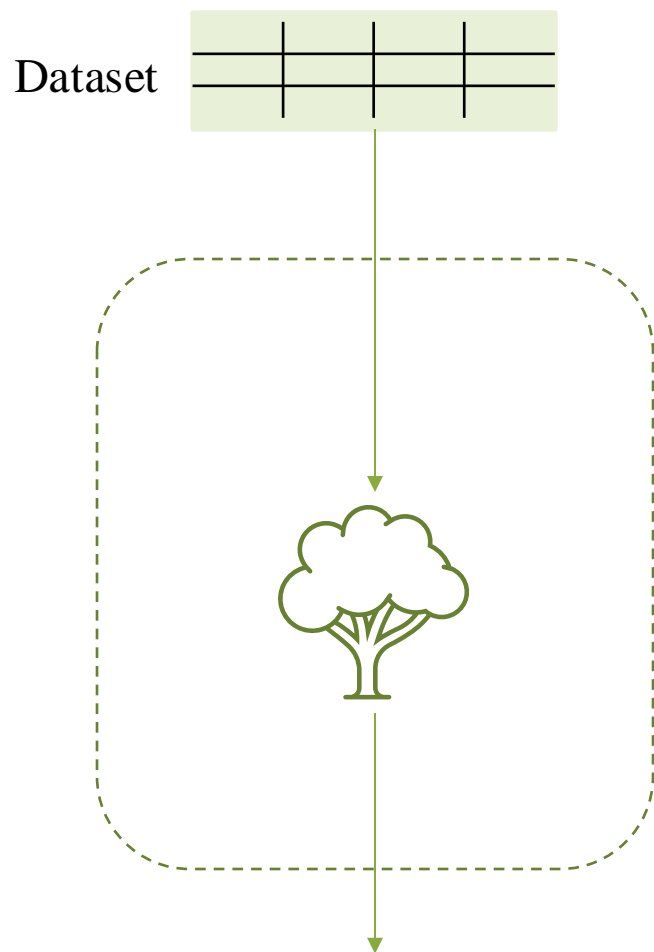


# Random Forest

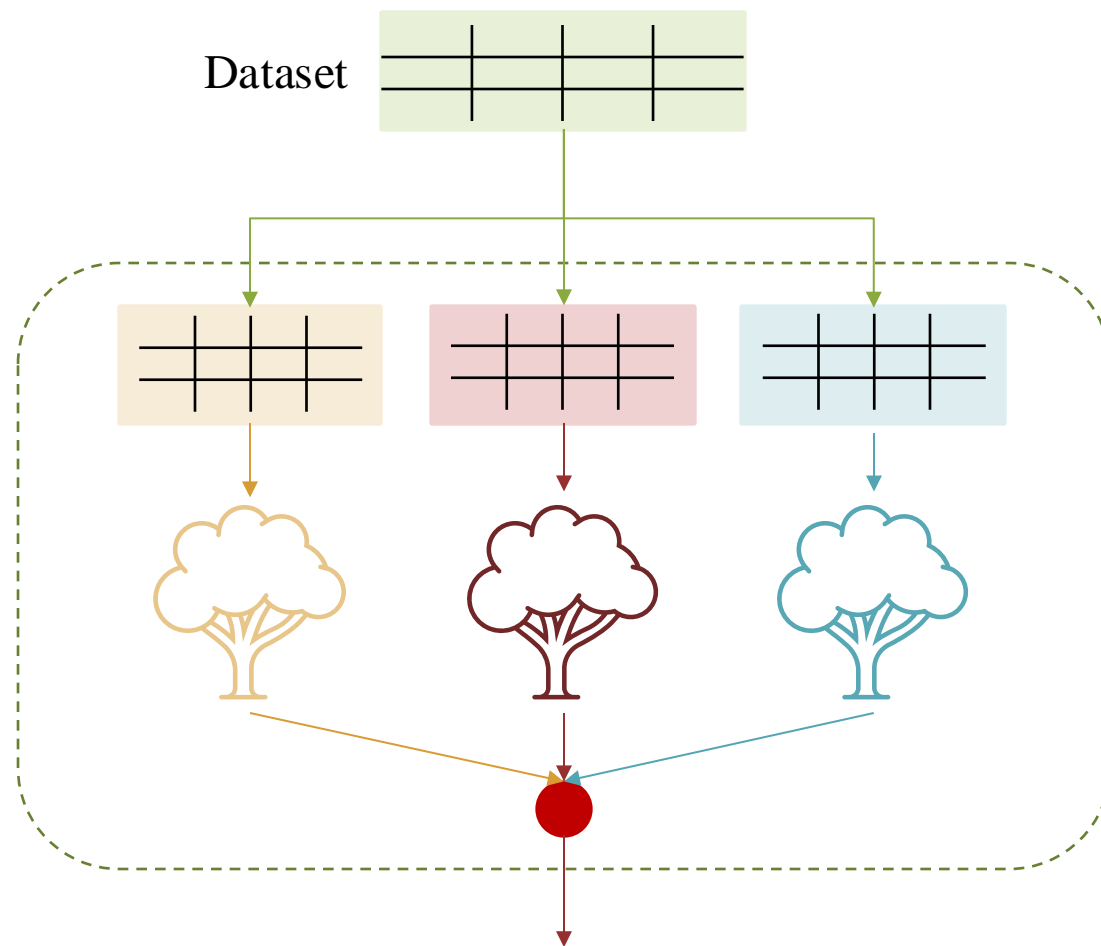


## From Decision Tree to Random Forest

Decision Tree



Random Forest

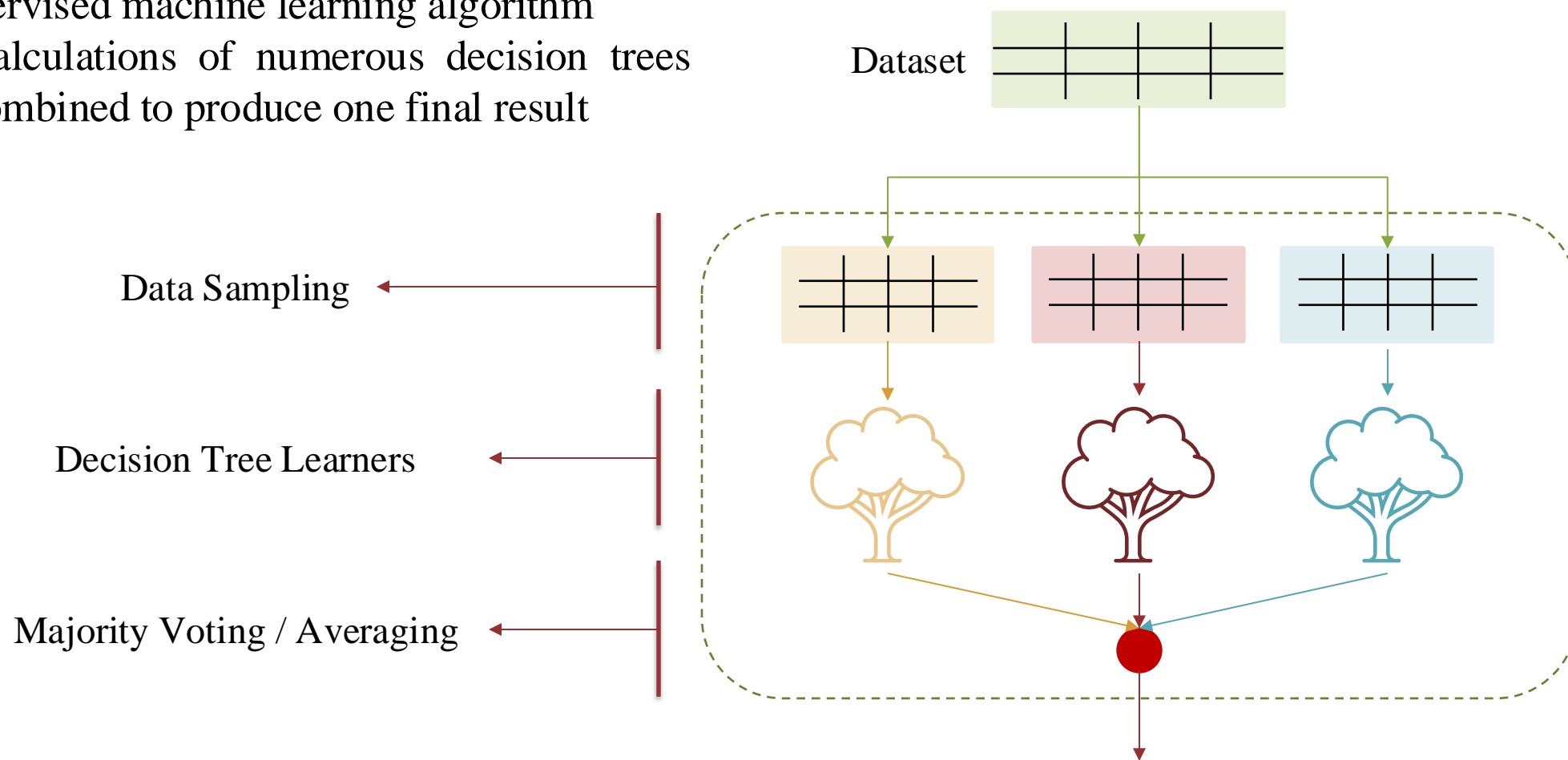


# Random Forest

## ! Random Forest

### A random forest

- ❖ a supervised machine learning algorithm
- ❖ the calculations of numerous decision trees are combined to produce one final result

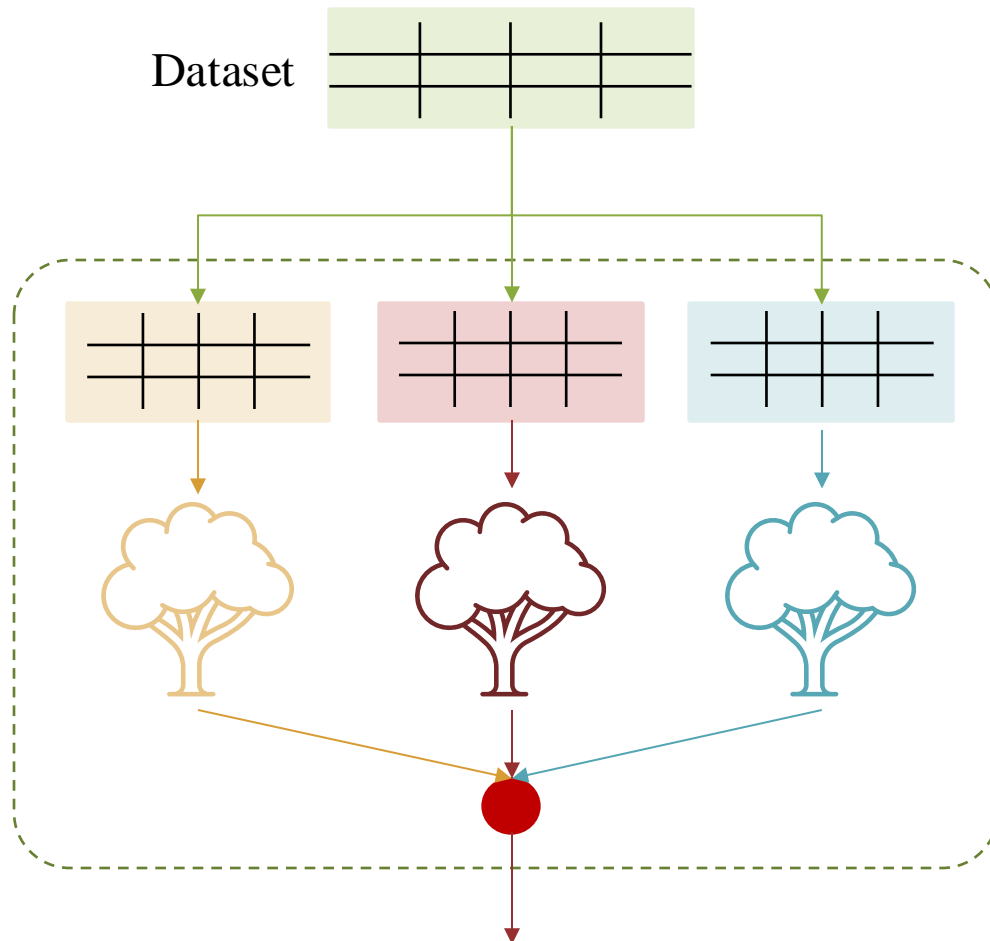




# Random Forest



## Data Sampling



Feature = 1  
Randomly sample  
with replacement

Length	Width	Label
1	0.2	0
1.3	0.6	0
0.9	0.7	0
1.7	0.5	1
1.8	0.9	1
1.2	1.3	1

Length	Label
1	0
1.3	0
1	0
1.8	1
1.8	1
1.2	1

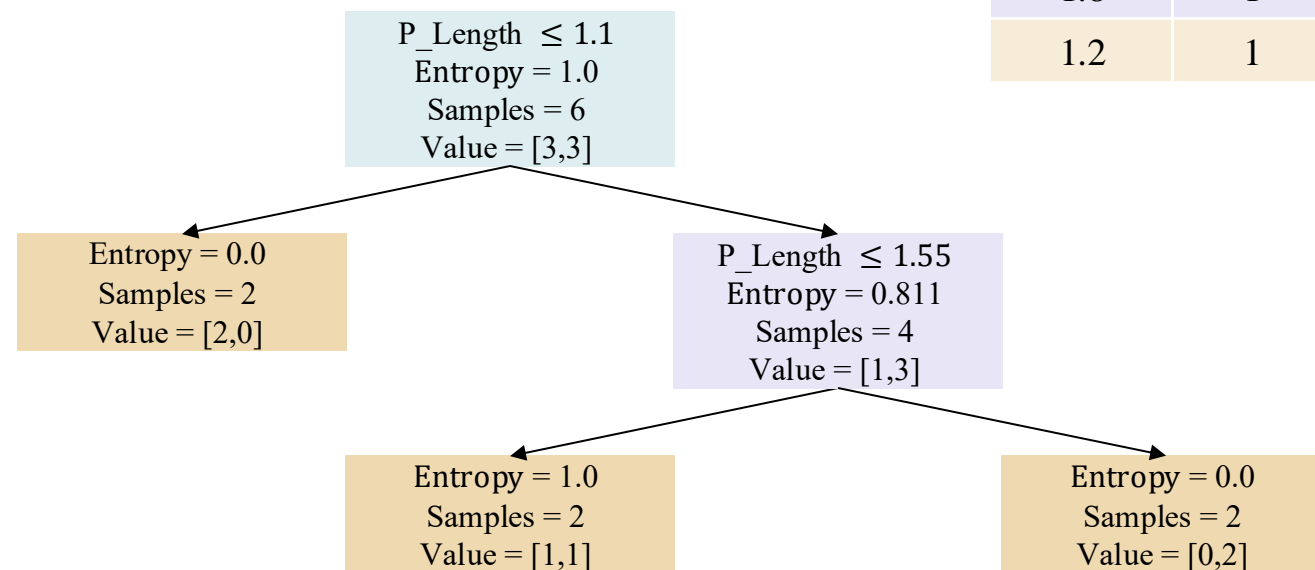
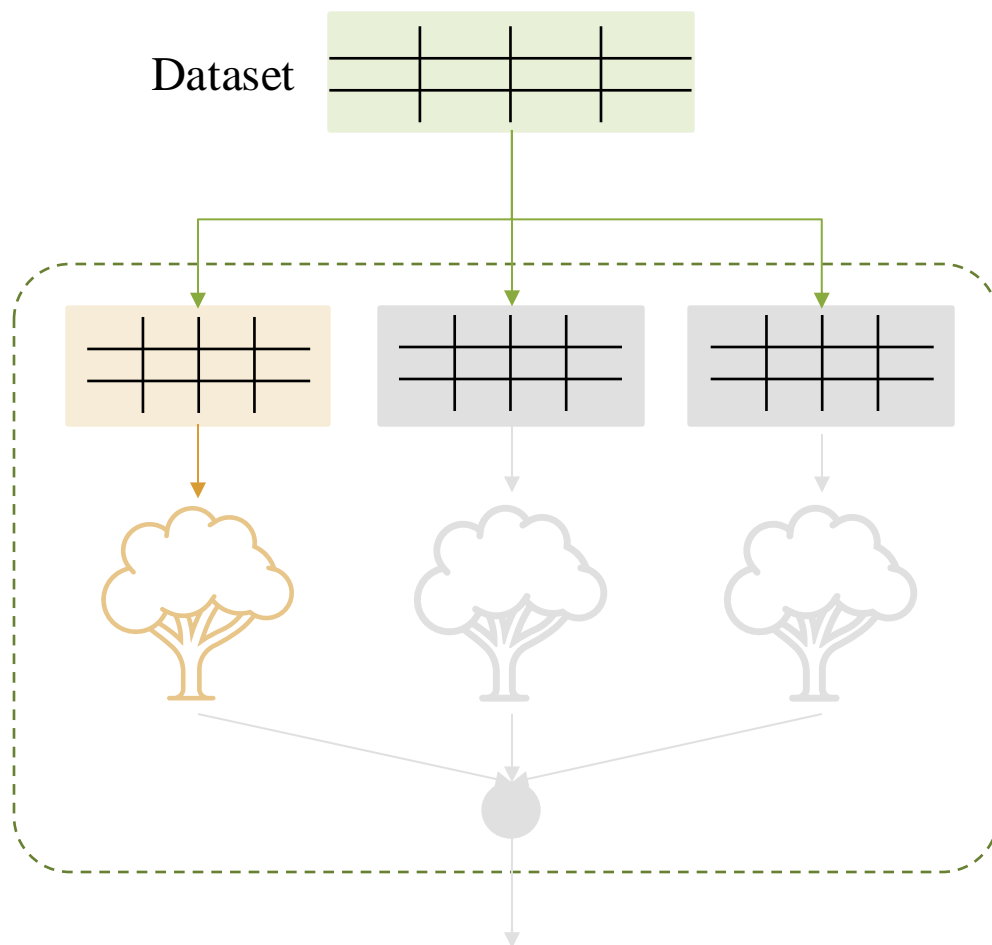
Width	Label
0.6	0
0.6	0
0.7	0
0.7	0
0.9	1
1.3	1

Length	Label
1	0
1.3	0
1.2	1
1.8	1
1.8	1
1.2	1

# Random Forest



## Decision Tree Learners

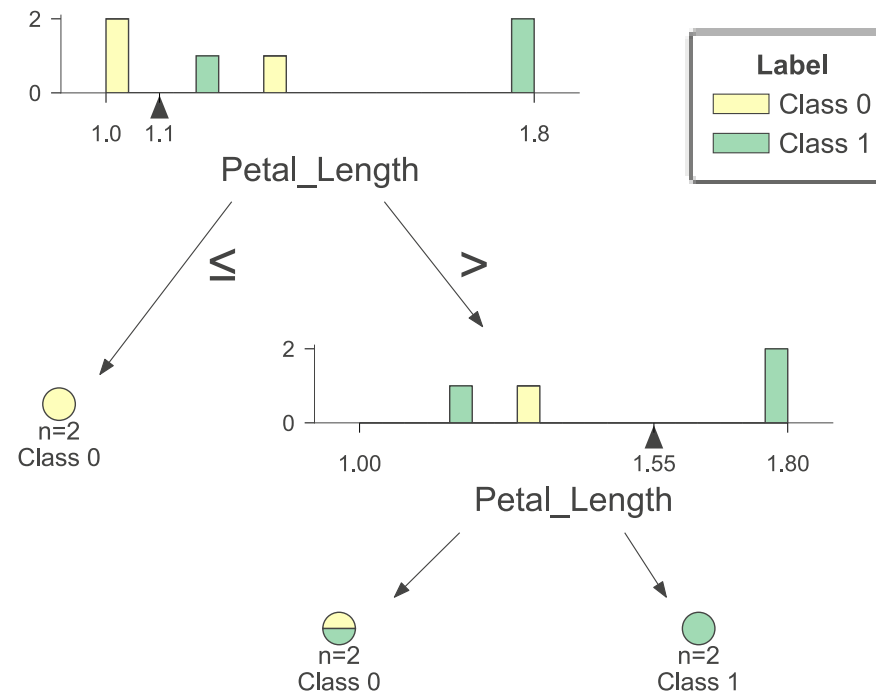
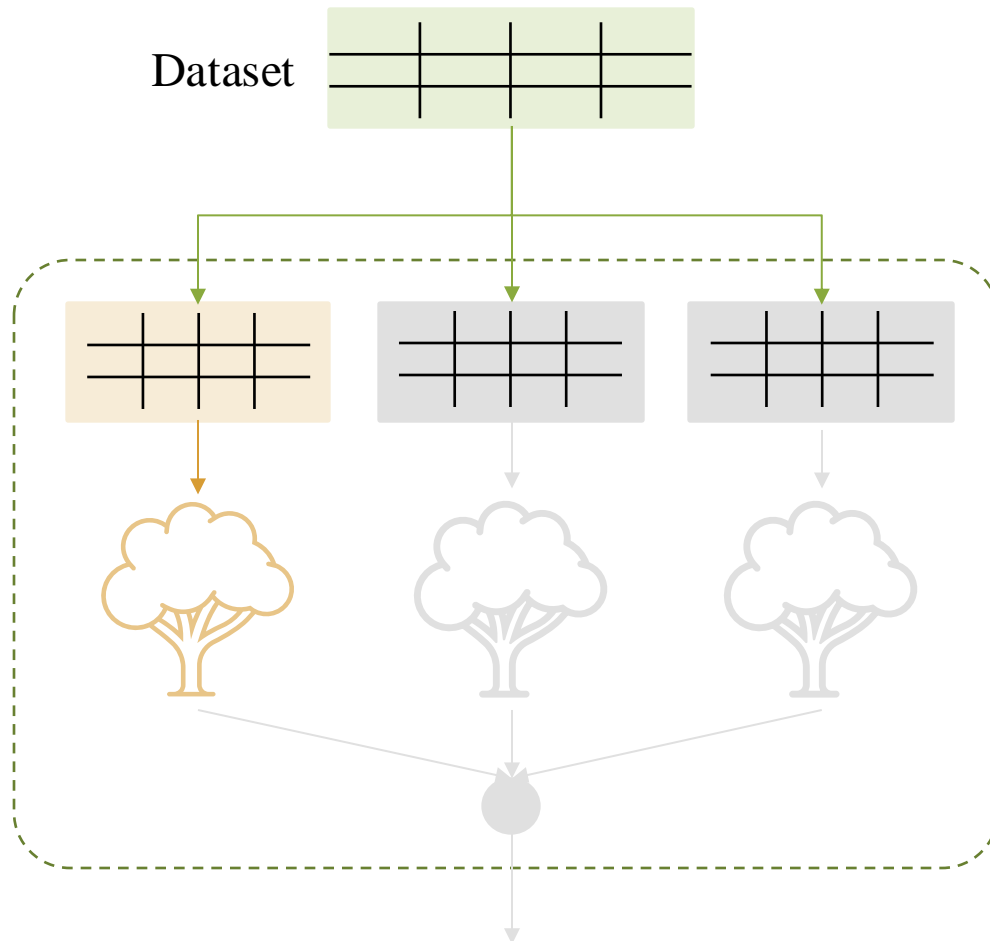


Length	Label
1	0
1.3	0
1	0
1.8	1
1.8	1
1.2	1

# Random Forest



## Decision Tree Learners

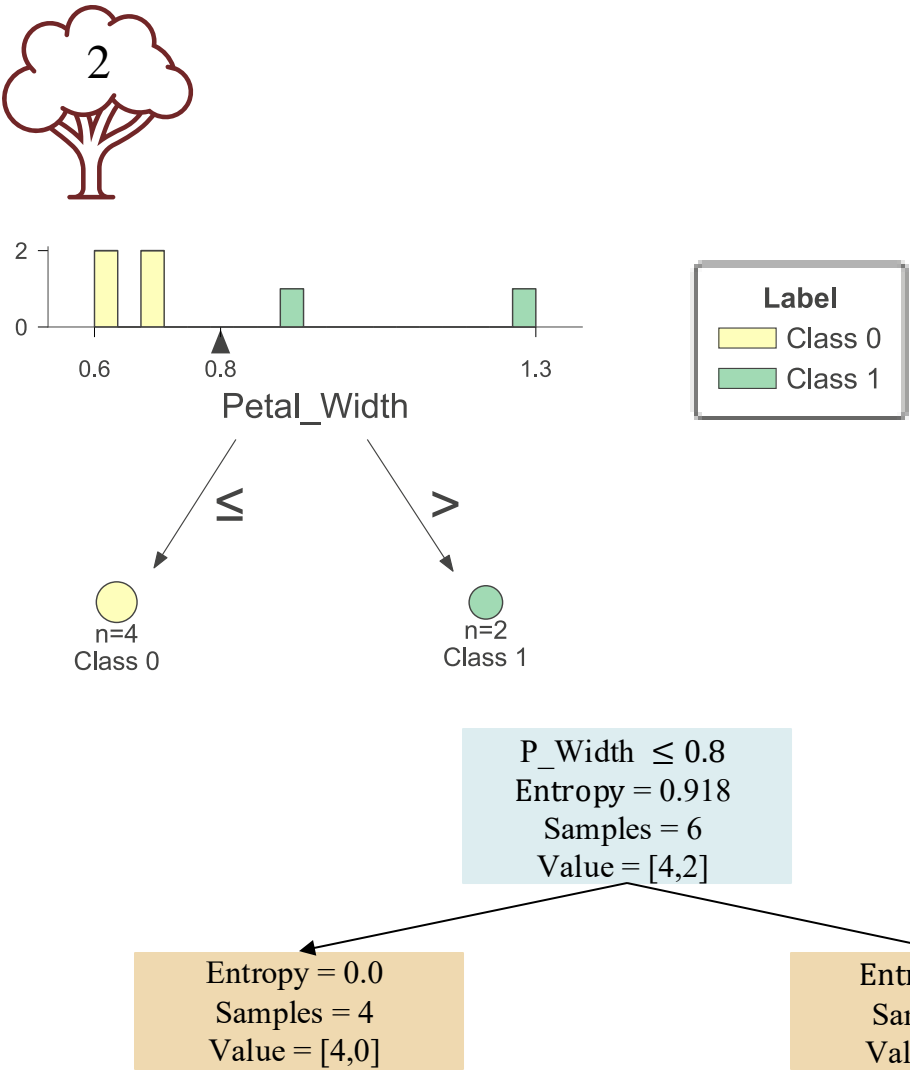
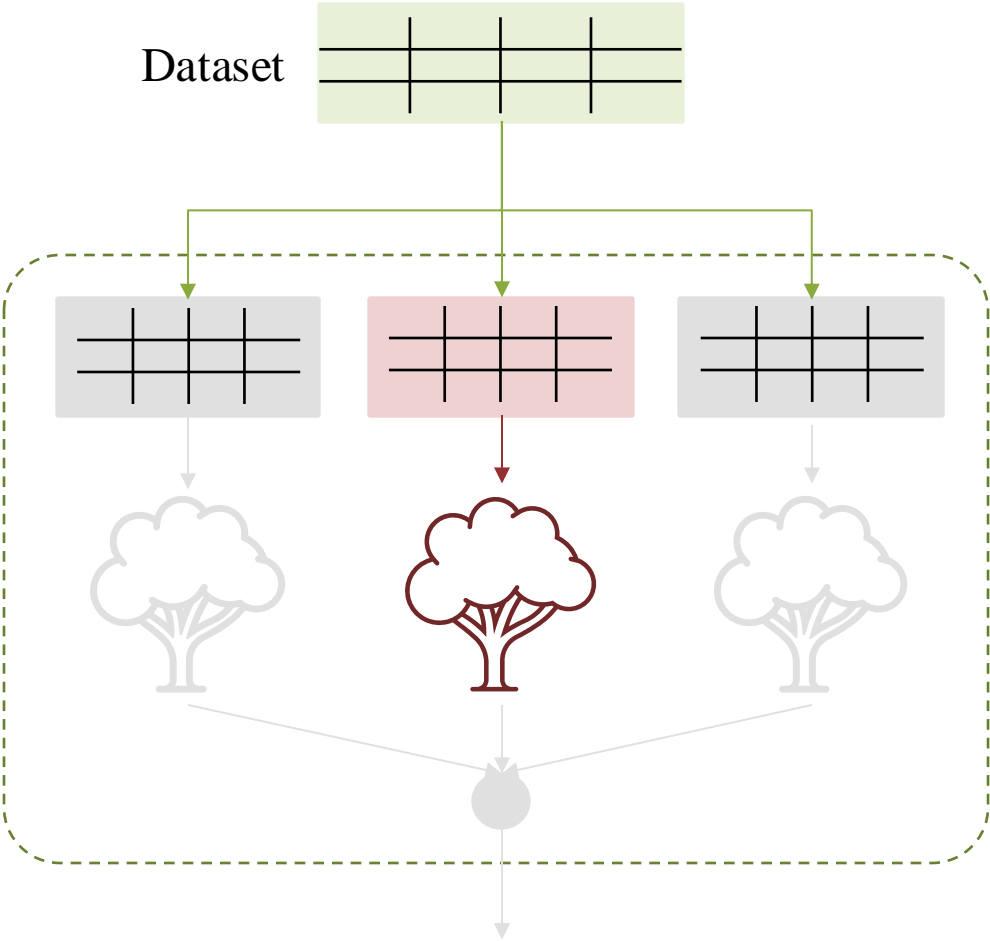


Length	Label
1	0
1.3	0
1	0
1.8	1
1.8	1
1.2	1

# Random Forest



## Decision Tree Learners



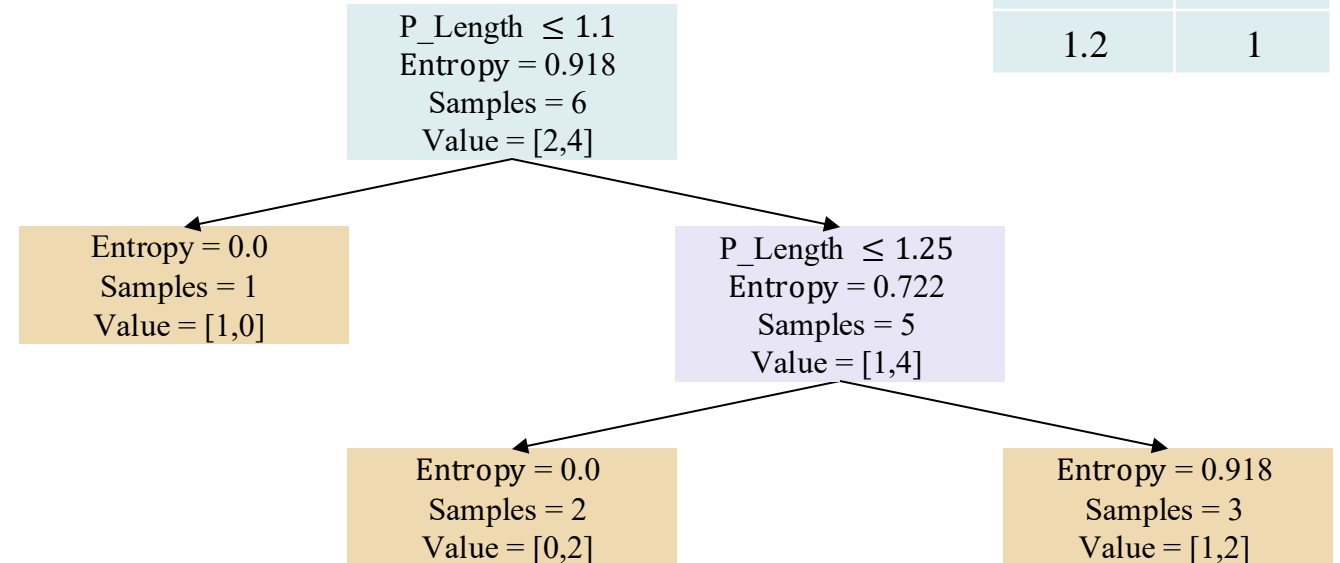
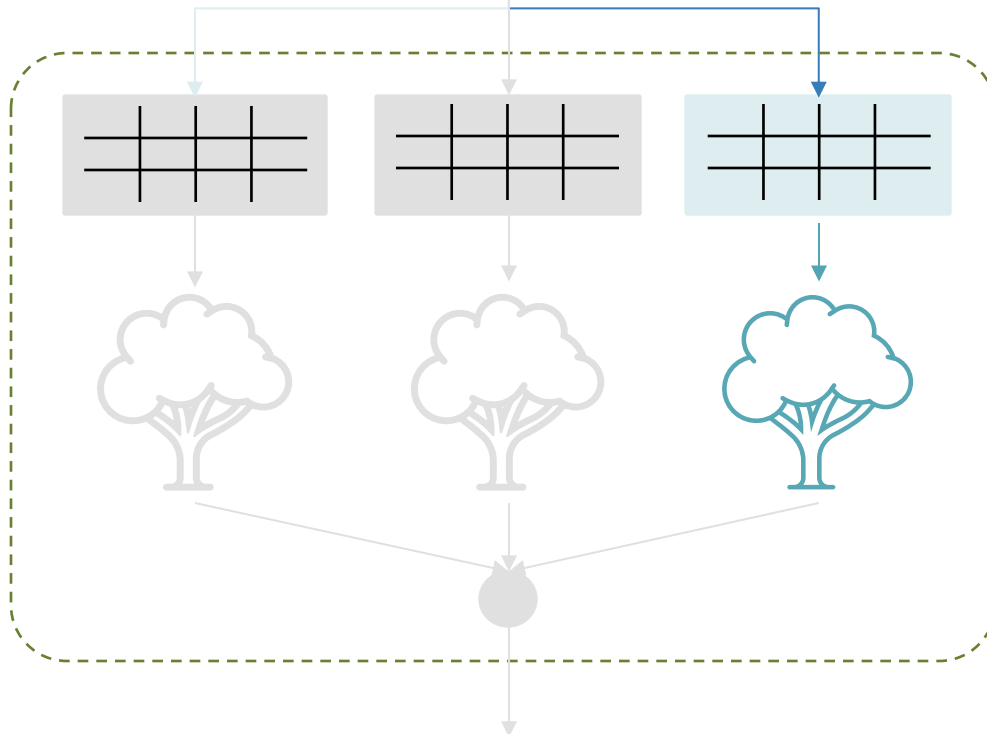
# Random Forest



## Decision Tree Learners



Dataset

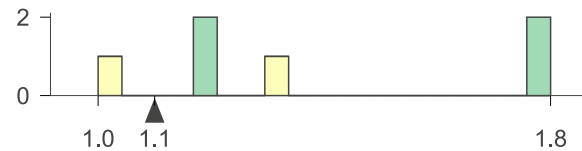
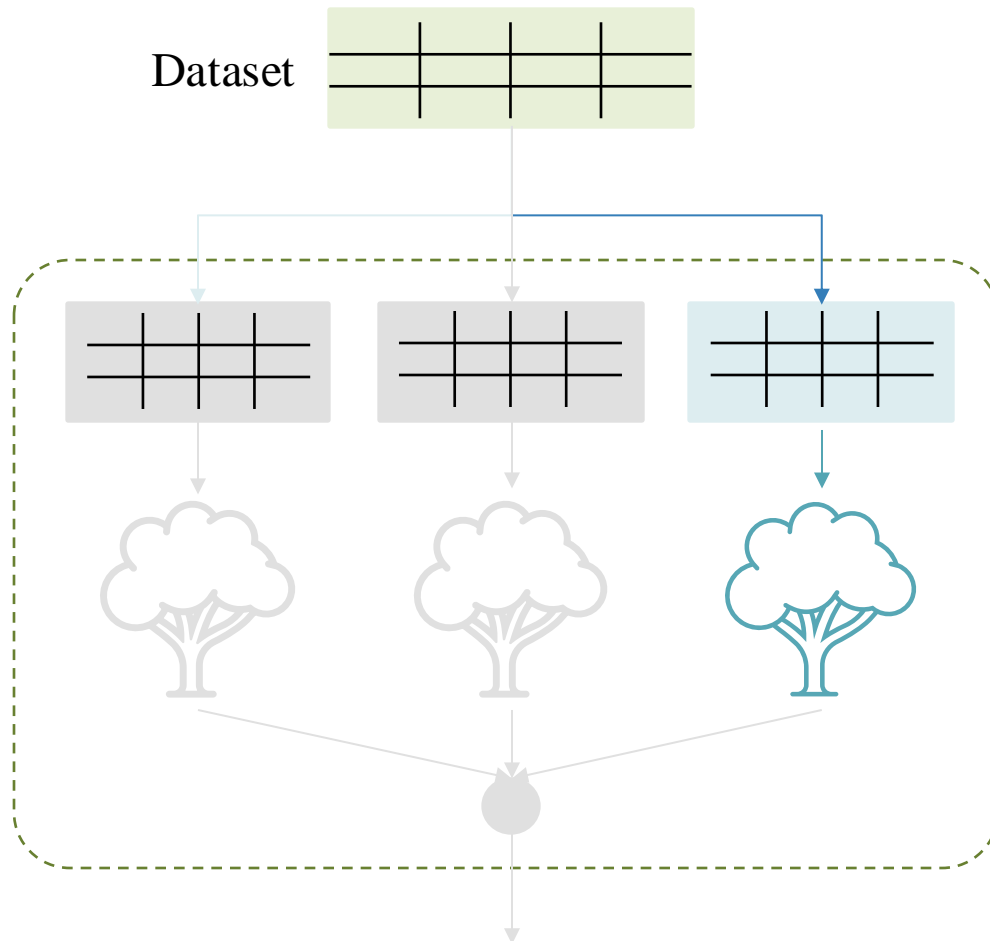



Length	Label
1	0
1.3	0
1.2	1
1.8	1
1.8	1
1.2	1

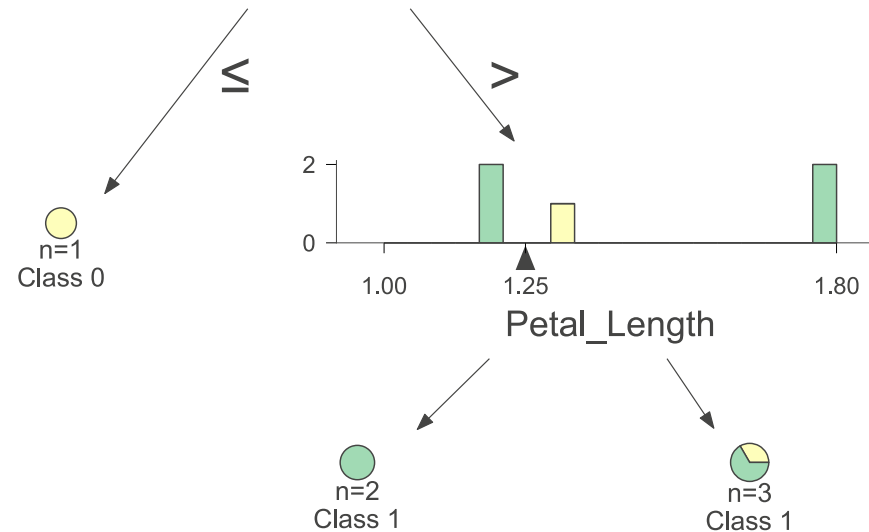
# Random Forest



## Decision Tree Learners



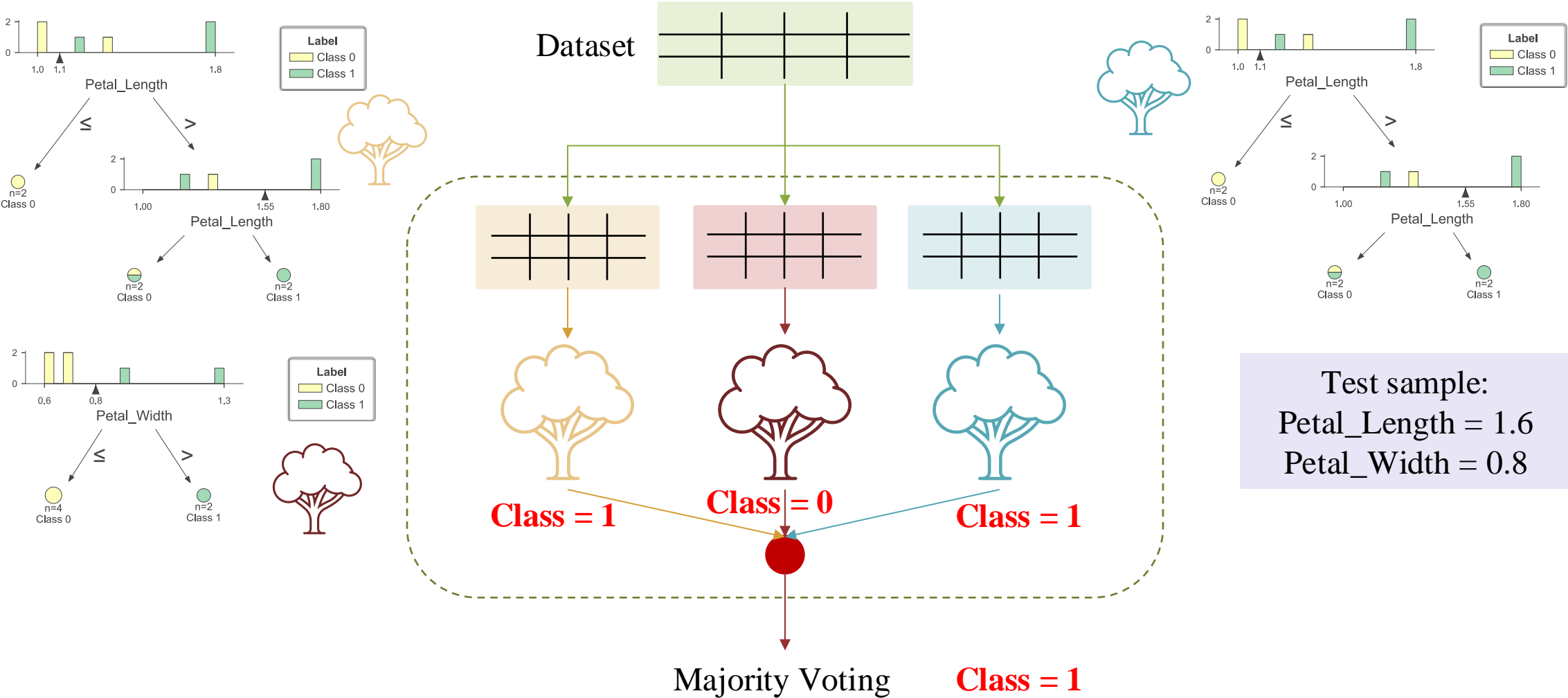
Length	Label
1	0
1.3	0
1.2	1
1.8	1
1.8	1
1.2	1



# Random Forest



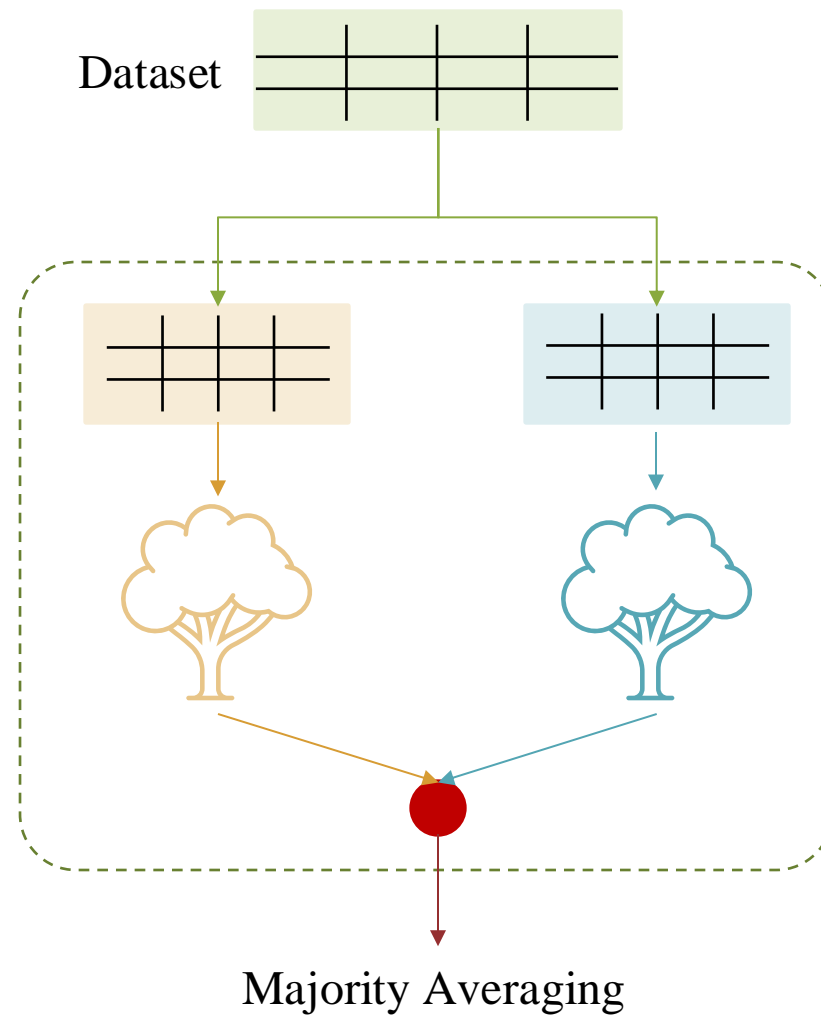
## Majority Voting for Classification



# Random Forest



## Random Forest for Regression

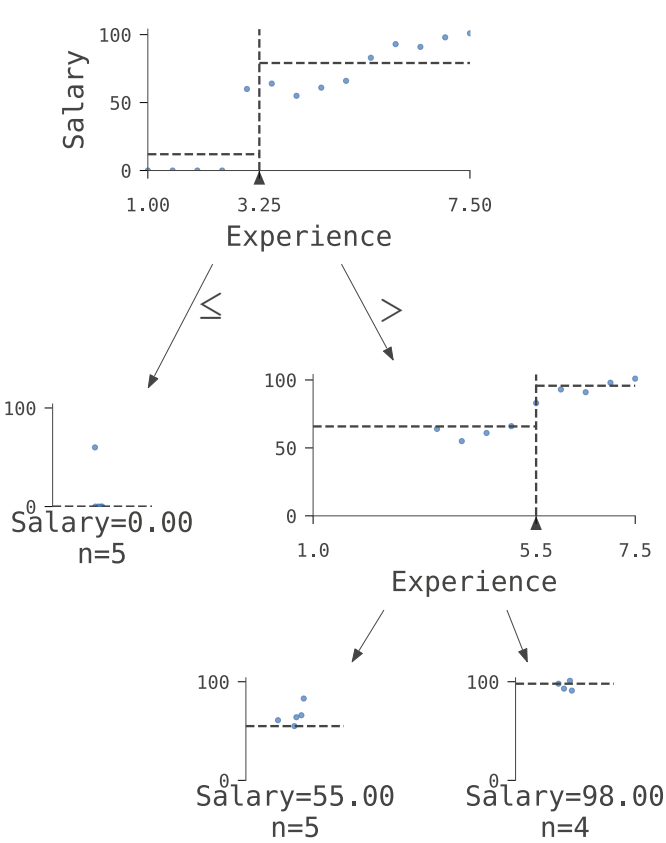




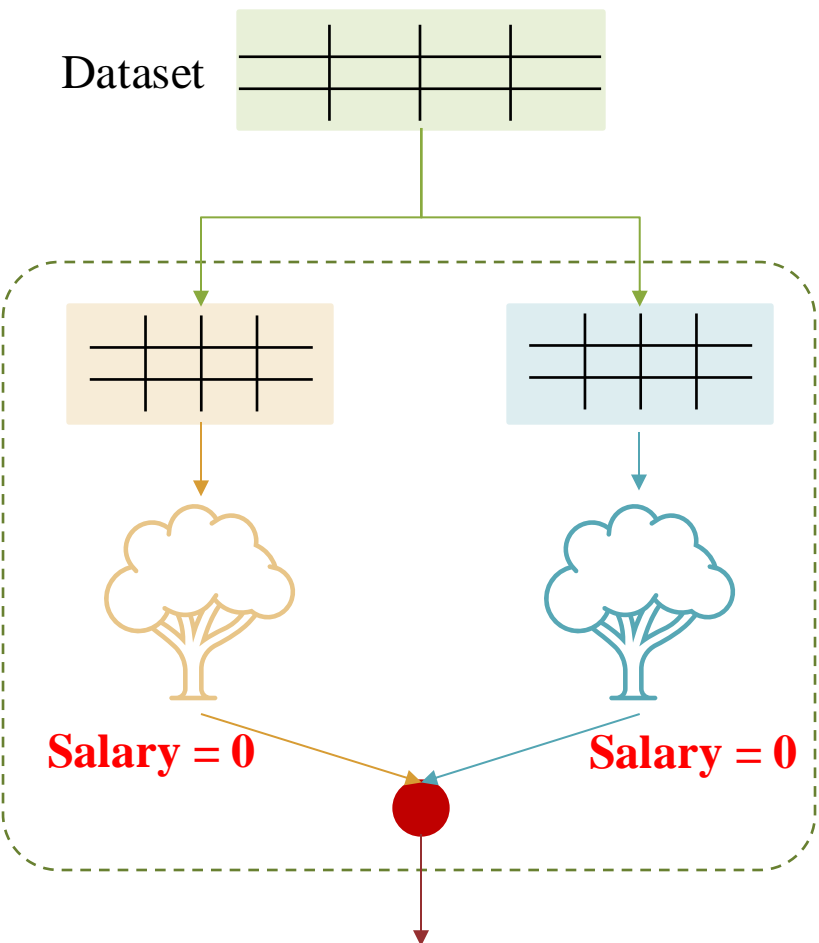
# Random Forest



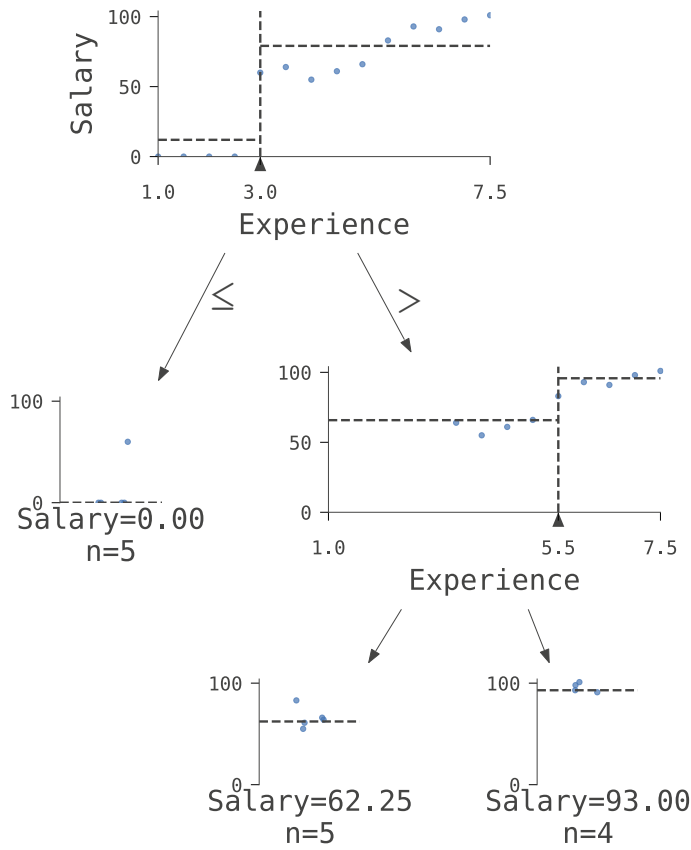
## Majority Averaging for Regression



Test sample:  
Experience = 3



Majority Averaging     **Salary = (0+0)/2=0**



QUIZ TIME

# Random Forest



## Random Forest for Classification

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.ensemble import RandomForestClassifier
4
5 # Load data
6 df = pd.read_csv('iris_2D.csv')
7
8 # Get data
9 x_data = df[['Petal_Length', 'Petal_Width']].to_numpy()
10 x_data = x_data.reshape(6, 2)
11 y_data = df['Label'].to_numpy(dtype=np.uint8)
```

```
13 # Define model
14 rf_classifier = RandomForestClassifier(
15     n_estimators=3,
16     max_features=1,
17     max_depth=1,
18     criterion='entropy',
19     max_samples=6
20 )
21
22 # Train model
23 rf_classifier.fit(x_data, y_data)
24
25 # Predict
26 x_test = np.array([[2.7, 0.8]])
27 y_predicted = rf_classifier.predict(x_test)
28 y_predicted
```

# Random Forest



## Random Forest for Regression

```
1 import pandas as pd
2 from sklearn.ensemble import RandomForestRegressor
3
4 # Load data
5 df = pd.read_csv('Salary_Data.csv')
6
7 # Get data
8 x_data = df.iloc[:, :-1]
9 y_data = df.iloc[:, -1]
10
```

```
11 # Define model
12 rf_regressor = RandomForestRegressor(
13     n_estimators=2,
14     max_depth=2,
15     max_samples=7
16 )
17
18 # Train model
19 rf_regressor.fit(x_data, y_data)
20
21 # Predict
22 x_test = x_data[:2]
23 y_pred = rf_regressor.predict(x_test)
24 y_pred
```

# Summary

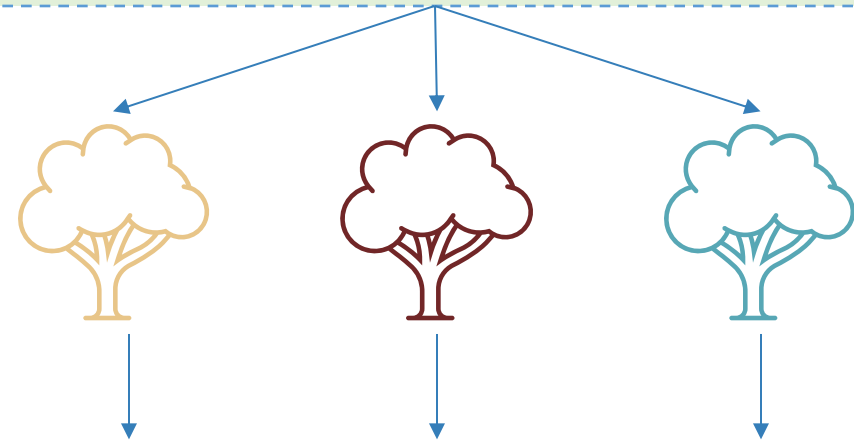
## Decision Tree Review

- ❖ Decision Tree
- ❖ Decision Tree for Classification
- ❖ Decision Tree for Regression
- ❖ IRIS Classification
- ❖ Salary Prediction



## Random Forest

- ❖ Decision Tree
- ❖ Random Forest
  - Bootstrap Sample
  - Majority Voting / Averaging
- ❖ IRIS Classification
- ❖ Salary Prediction





AI VIET NAM

@aivietnam.edu.vn

# Thanks!

## Any questions?