

---

# CS 484

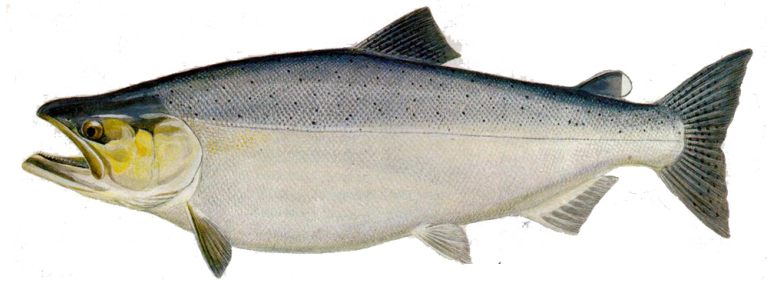
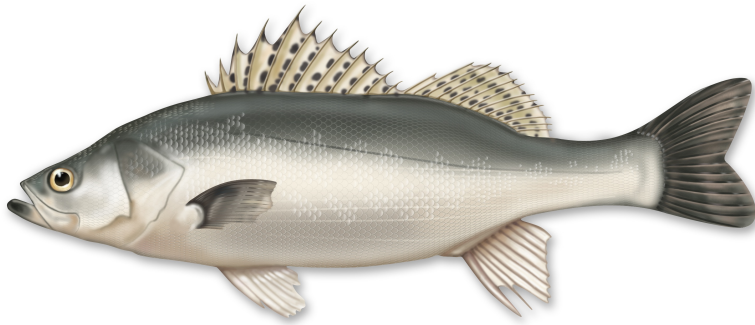
# Data Mining

Classification 2

---

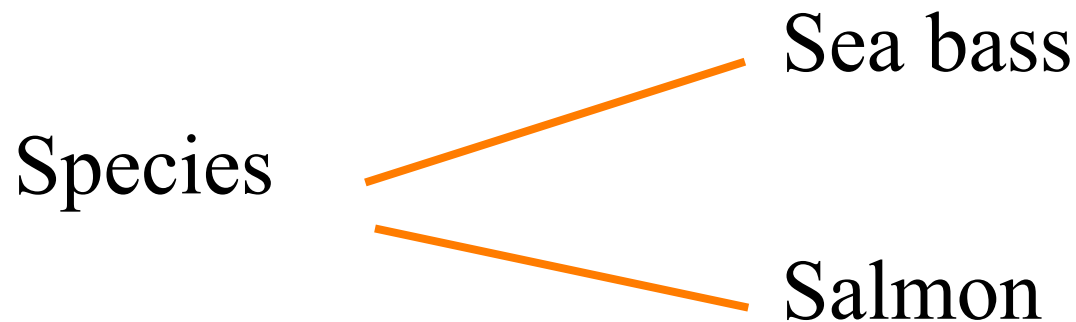
---

# Case Study: Sea Bass vs Salmon ?



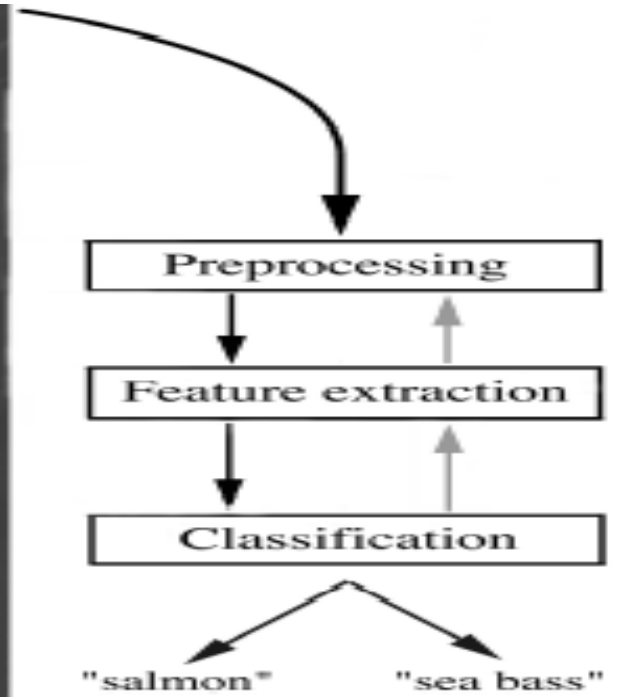
# An Example

- “Sorting incoming Fish on a conveyor according to species using optical sensing”



# Problem Analysis

- Set up a camera and take some sample images to extract features
  - Length
  - Lightness
  - Width
  - Number and shape of fins
  - Position of the mouth, etc...

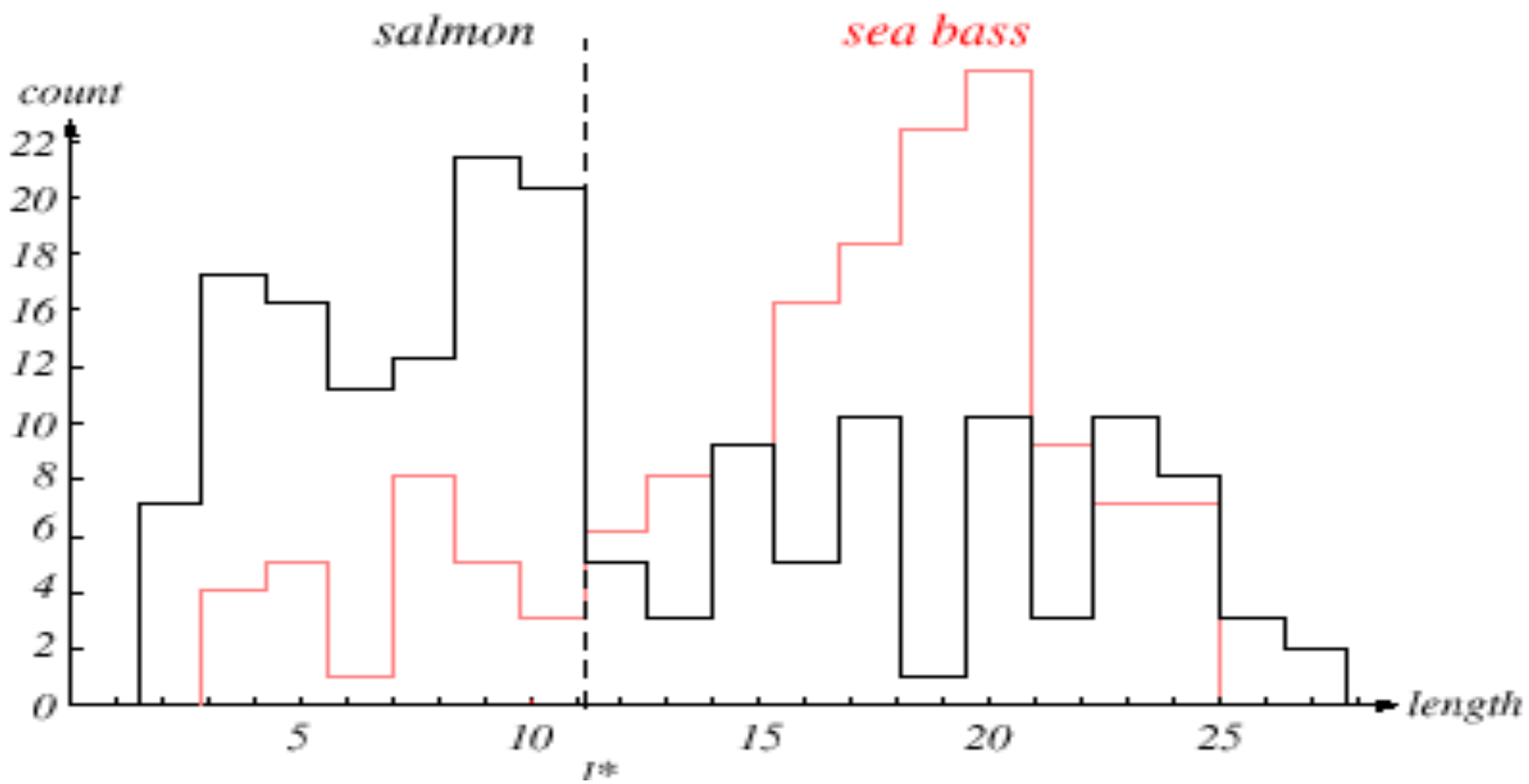


---

- Classification

- Select the length of the fish as a possible feature for discrimination





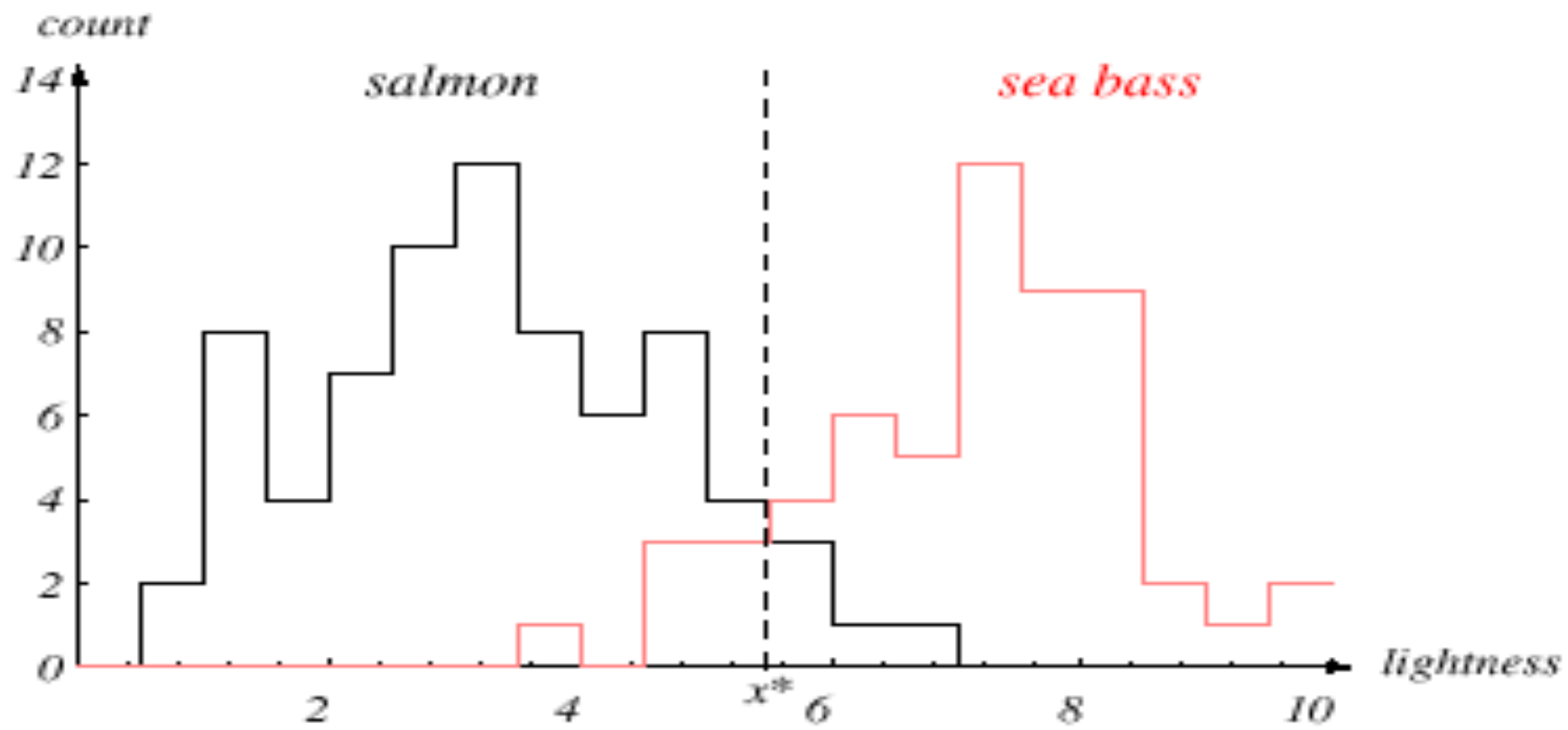


The **length** is a poor feature alone!

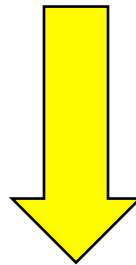
Select the **lightness** as a possible feature.







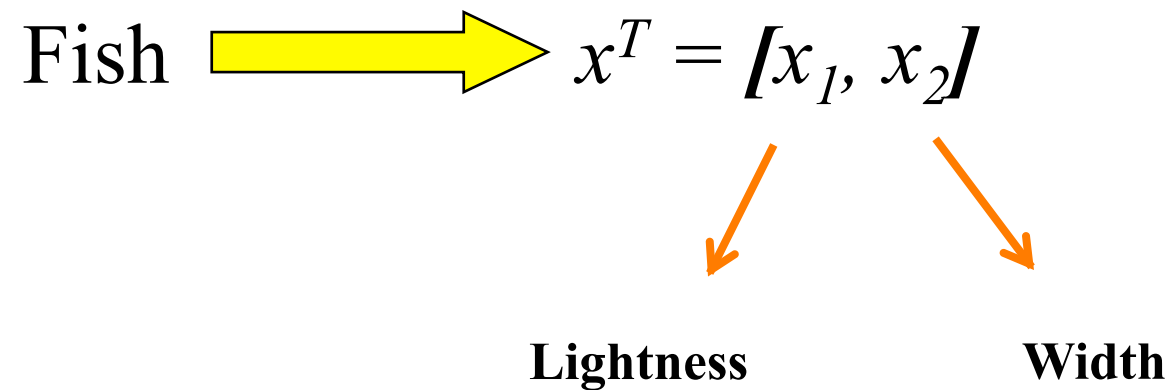
- 
- Threshold decision boundary and cost relationship
    - Move our decision boundary toward smaller values of lightness in order to minimize the cost (reduce the number of sea bass that are classified salmon!)

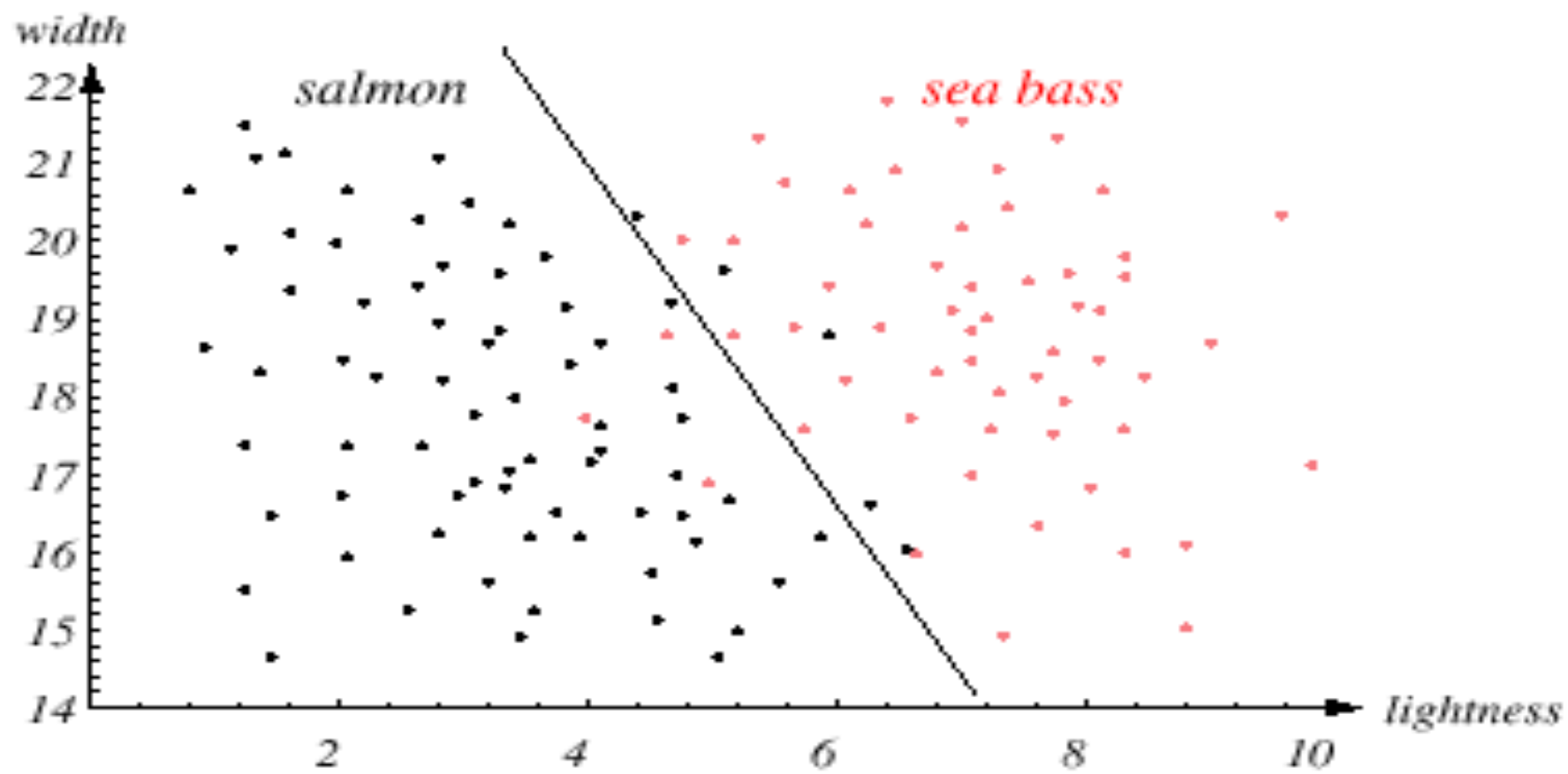




Task of decision theory

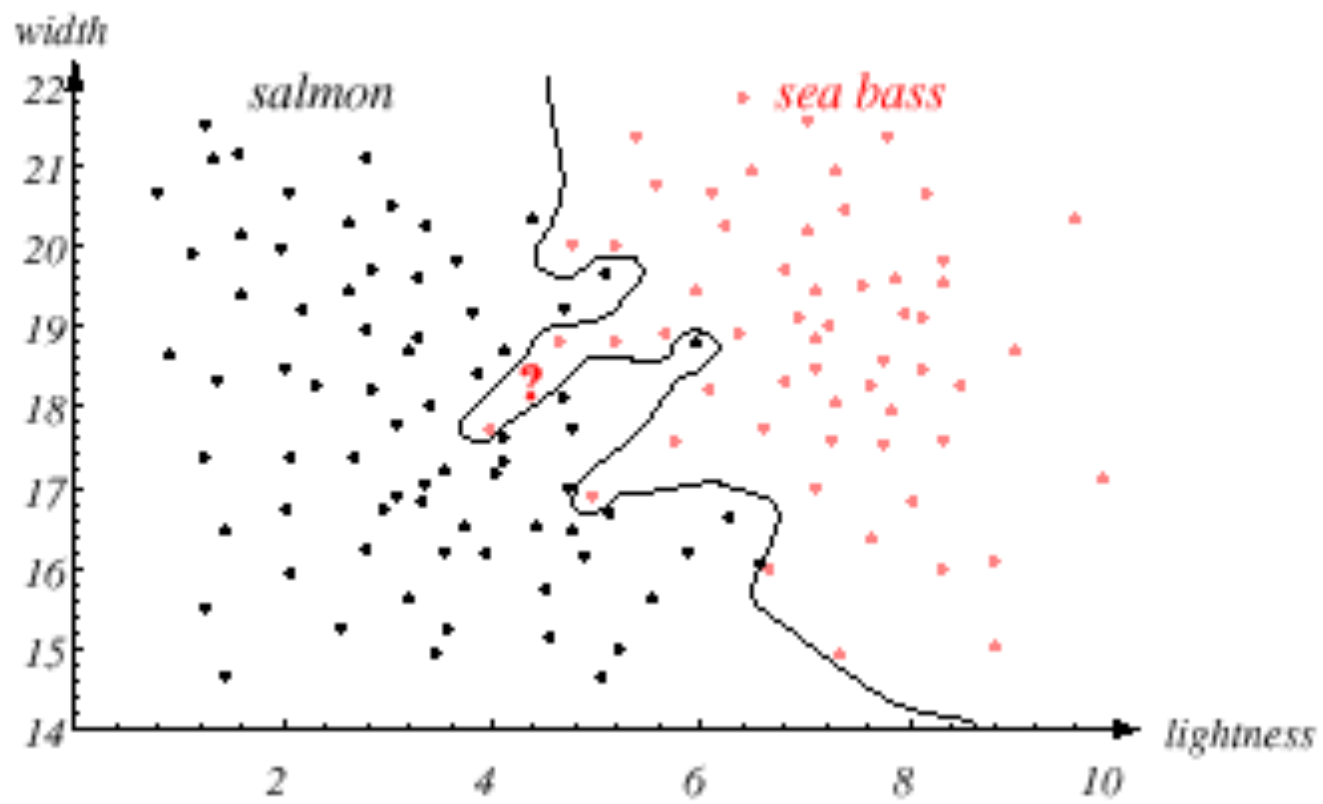
---

- 
- Adopt the lightness and add the width of the fish





- 
- We might add other features that are not correlated with the ones we already have. A precaution should be taken not to reduce the performance by adding such “noisy features”
  - Ideally, the best decision boundary should be the one which provides an optimal performance such as in the following figure:
- 

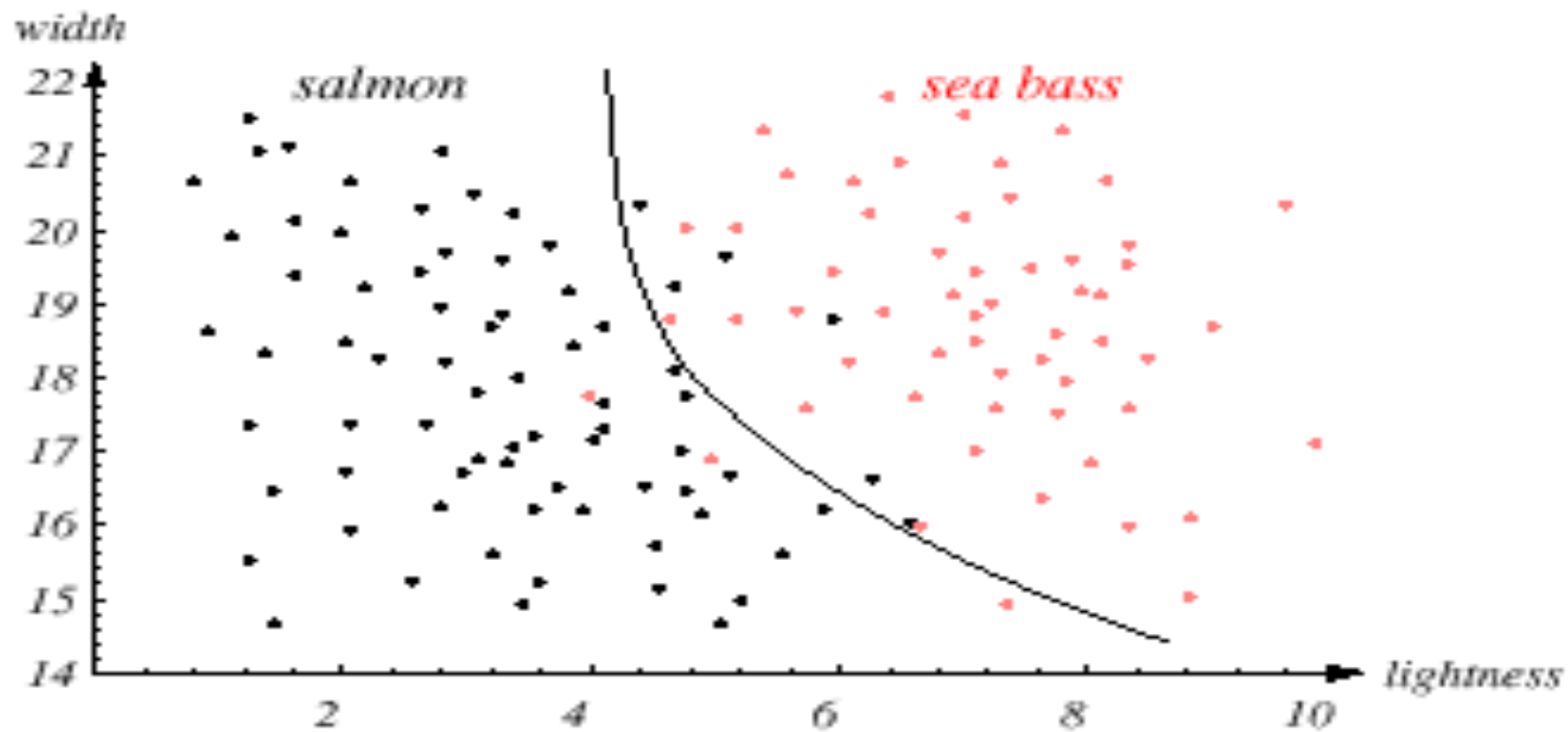


- 
- However, our satisfaction is premature because the central aim of designing a classifier is to correctly classify novel input



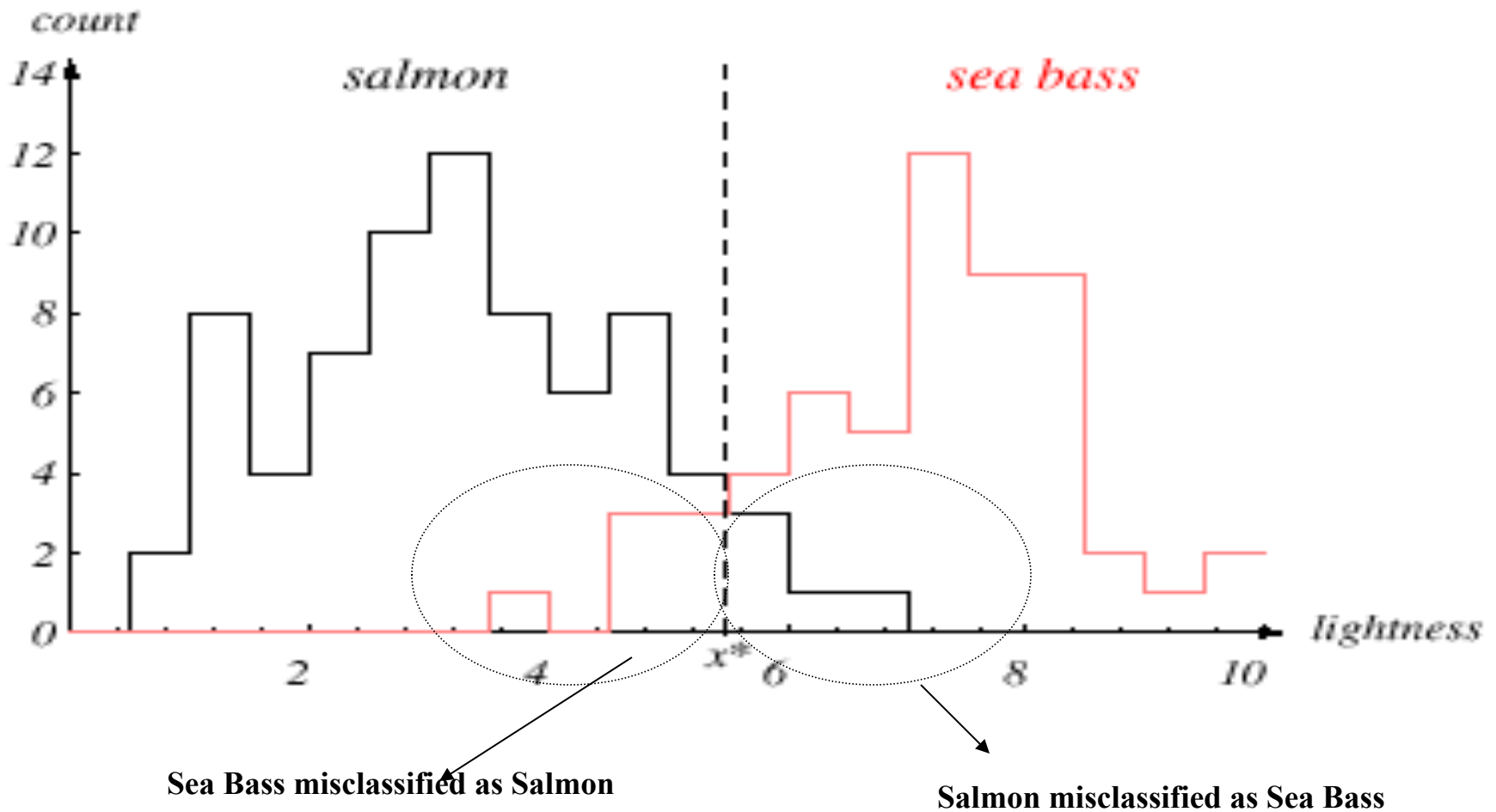
Issue of generalization!

---





# Misclassifications



---

# Cost sensitive classification

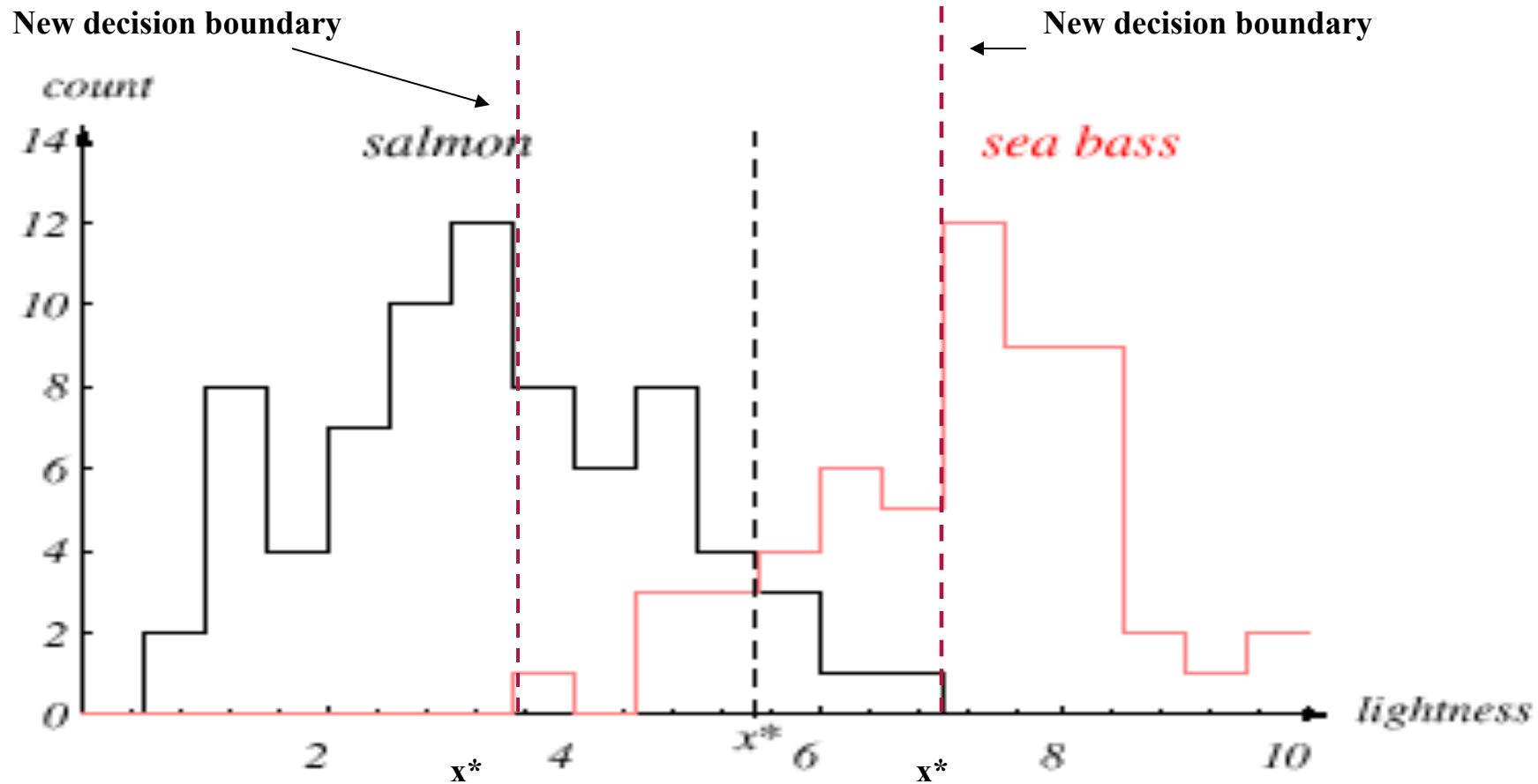
- Penalize misclassifications of one class more than the other
- Changes decision boundaries



**What if?**


Salmon is more expensive than Bass?

Bass is more expensive than Salmon?





# Classification Techniques

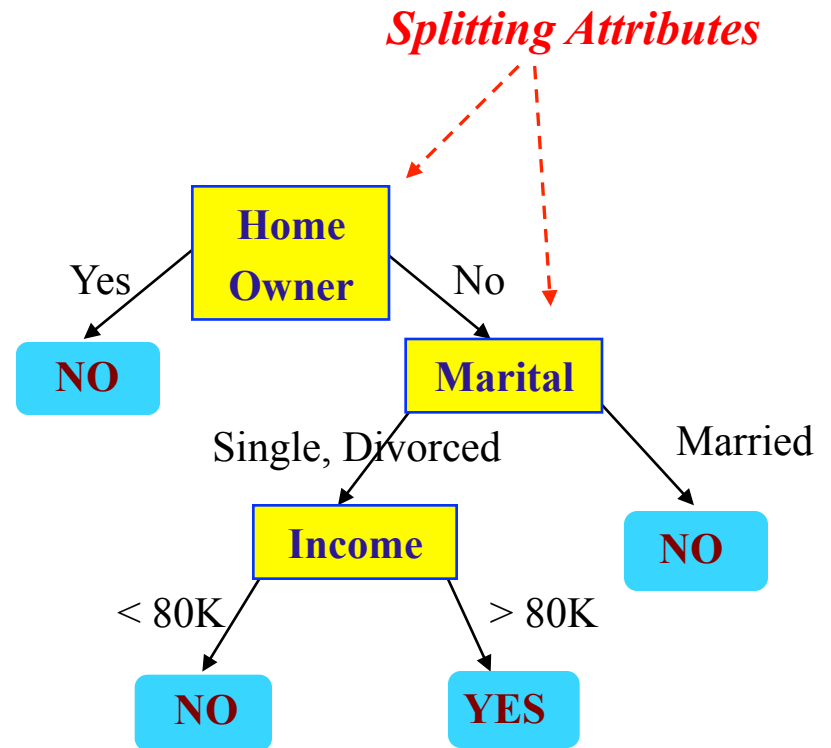
- Decision Tree based Methods
  - Rule-based Methods
  - Memory based reasoning
  - Neural Networks
  - Naïve Bayes and Bayesian Belief Networks
  - Support Vector Machines
- 

# Example of a Decision Tree

*categorical*      *categorical*      *continuous*      *class*

<i>Tid</i>	Home Owner	Marital Status	Annual Income	Defaulted
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

**Training Data**

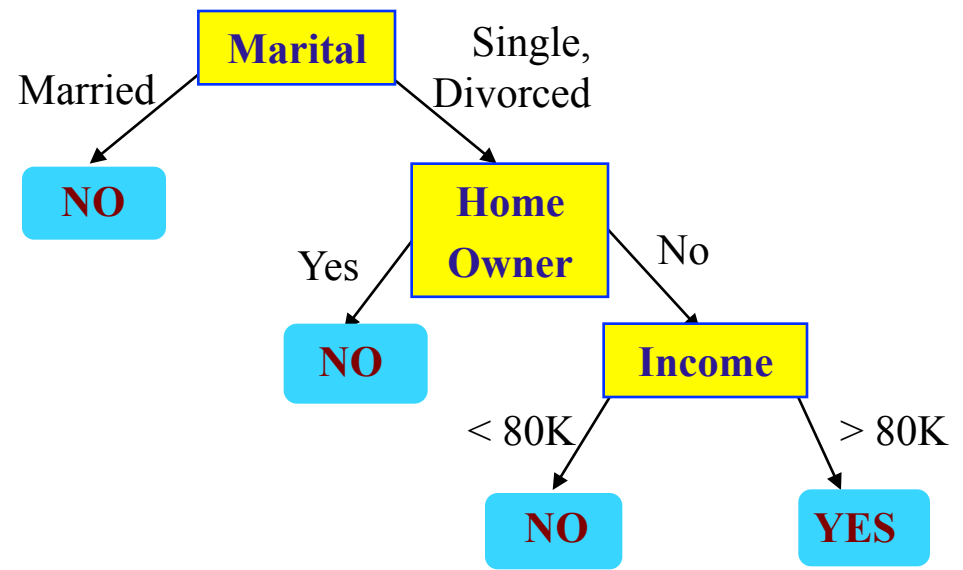


**Model: Decision Tree**

# Another Example of Decision Tree

*categorical*      *categorical*      *continuous*      *class*

<i>Tid</i>	Home Owner	Marital Status	Annual Income	Defaulted
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



**There could be more than one tree that fits the same data!**

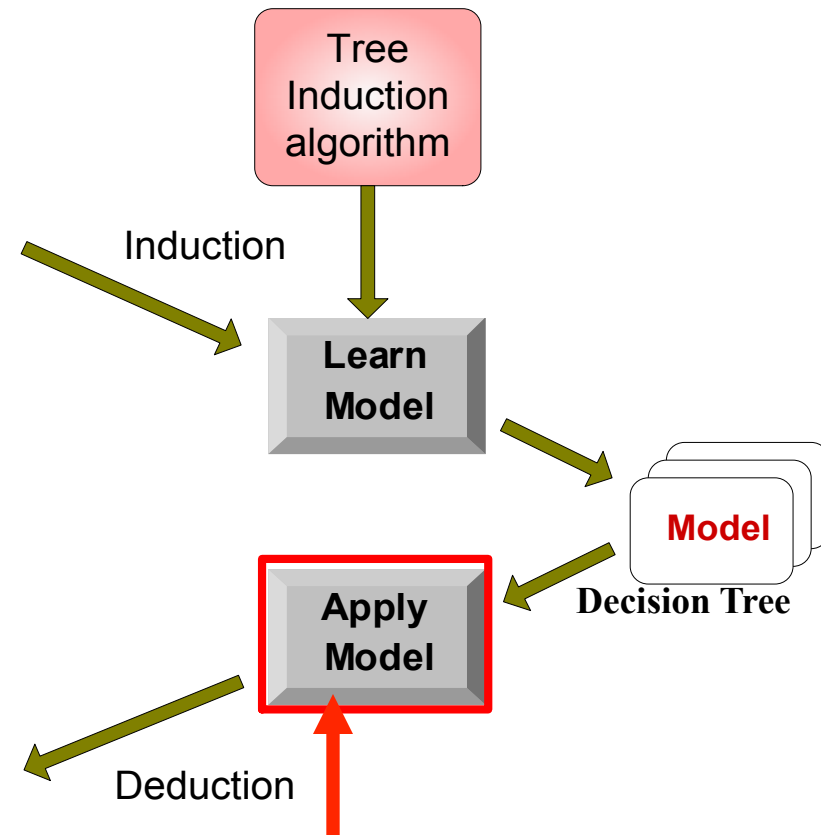
# Decision Tree Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set

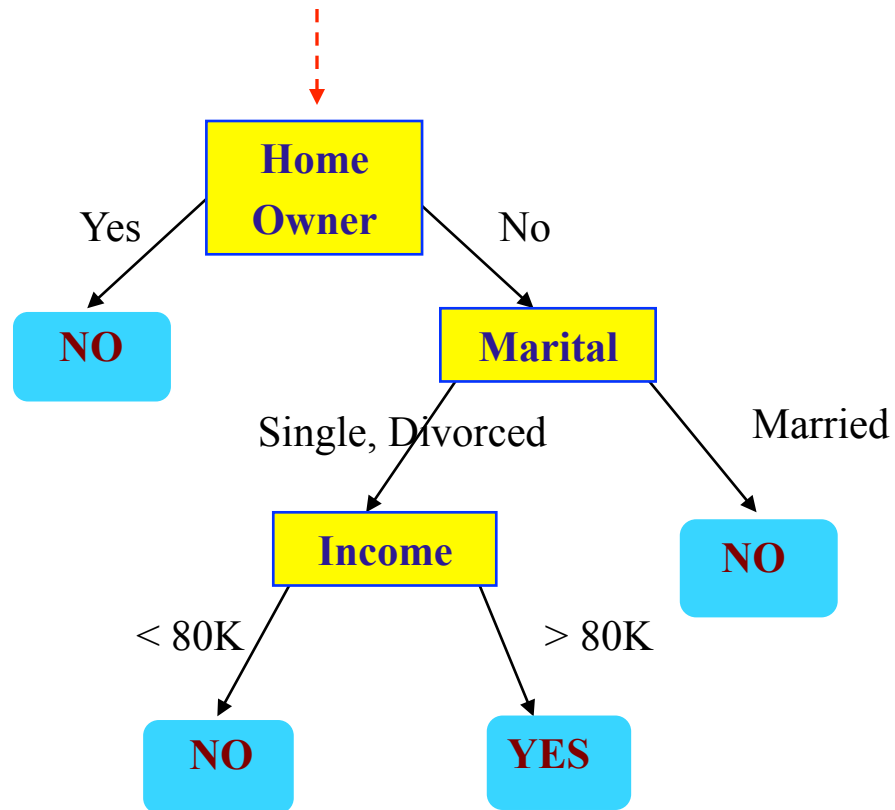


# Apply Model to Test Data

## Test Data

Home Owner	Marital Status	Annual Income	Defaulted
No	Married	80K	?

Start from the root of tree.

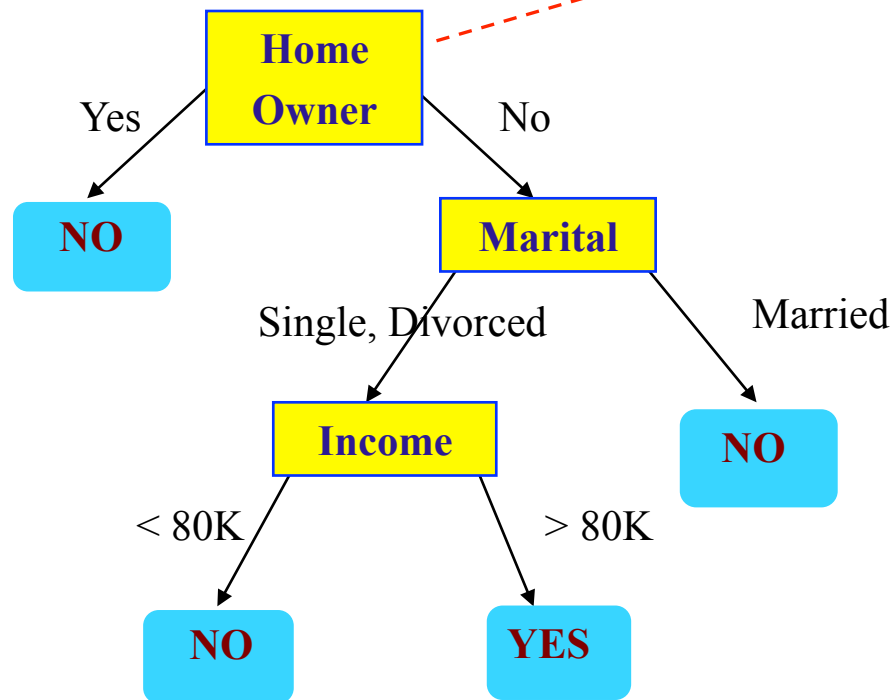




# Apply Model to Test Data

Test Data

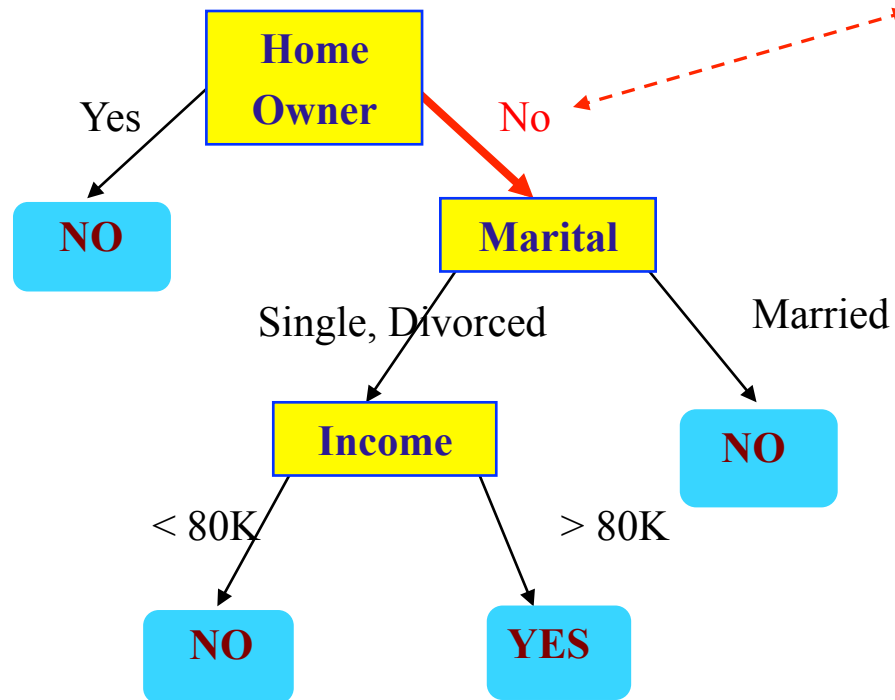
Home Owner	Marital Status	Annual Income	Defaulted
No	Married	80K	?



# Apply Model to Test Data

Test Data

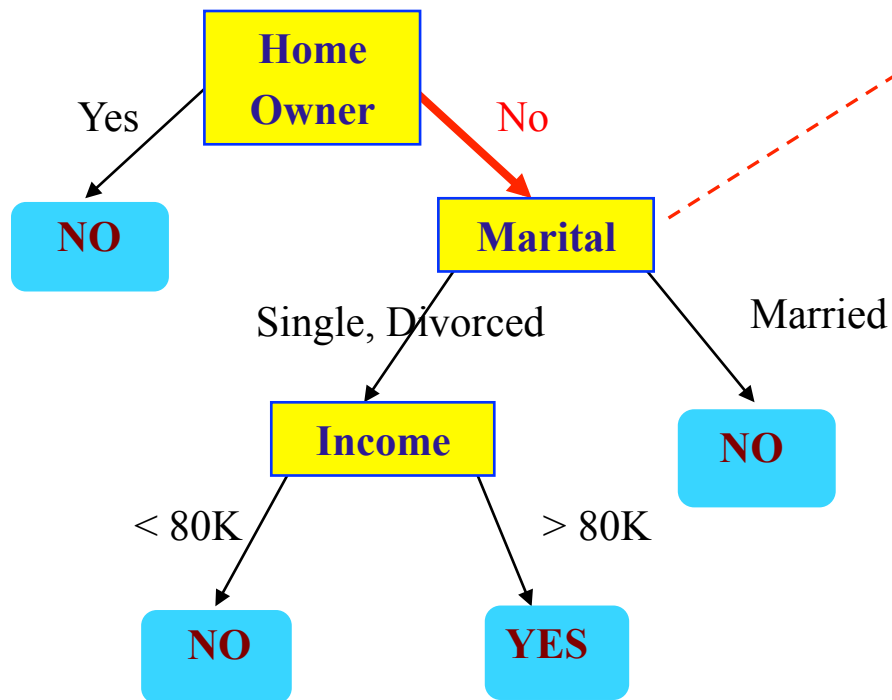
Home Owner	Marital Status	Annual Income	Defaulted
No	Married	80K	?



# Apply Model to Test Data

Test Data

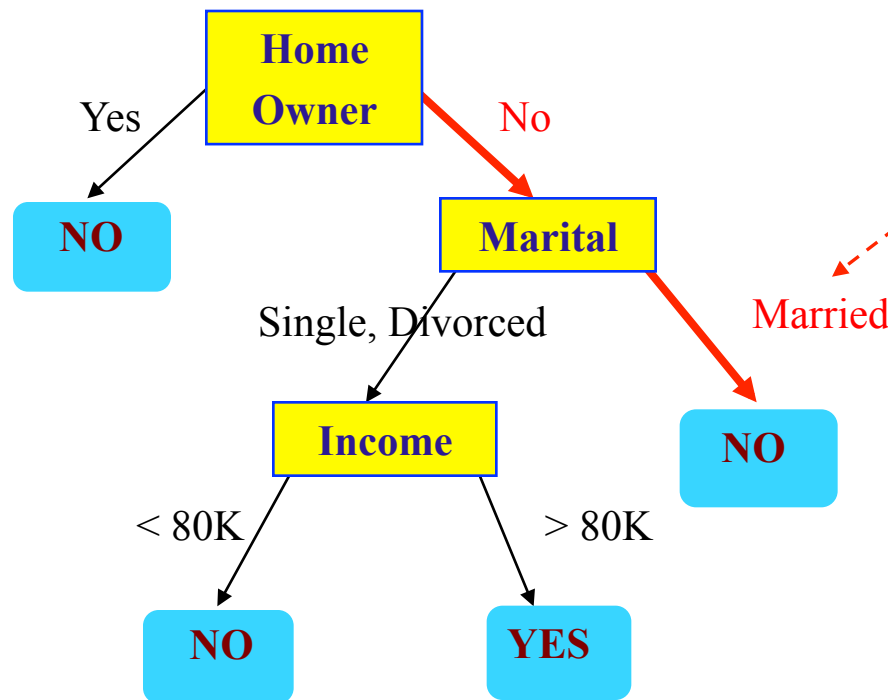
Home Owner	Marital Status	Annual Income	Defaulted
No	Married	80K	?



# Apply Model to Test Data

Test Data

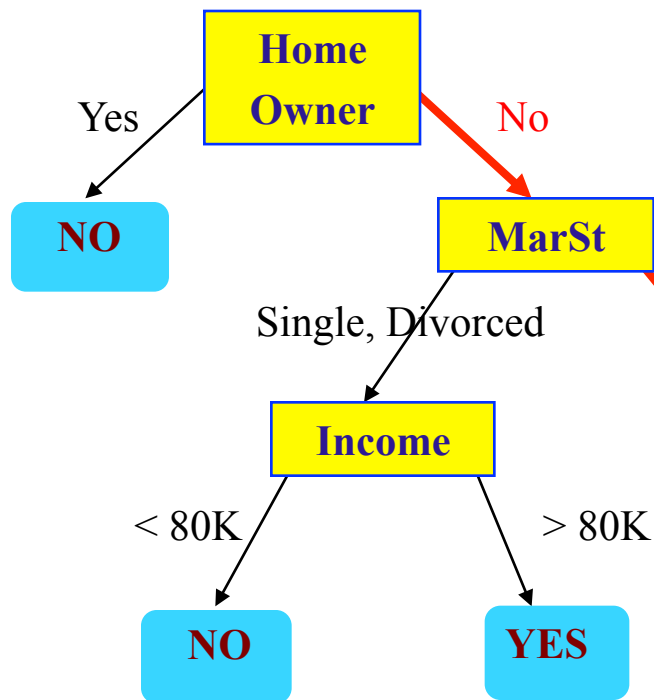
Home Owner	Marital Status	Annual Income	Defaulted
No	Married	80K	?



# Apply Model to Test Data

Test Data

Home Owner	Marital Status	Annual Income	Defaulted
No	Married	80K	?



Assign Defaulted to "No"

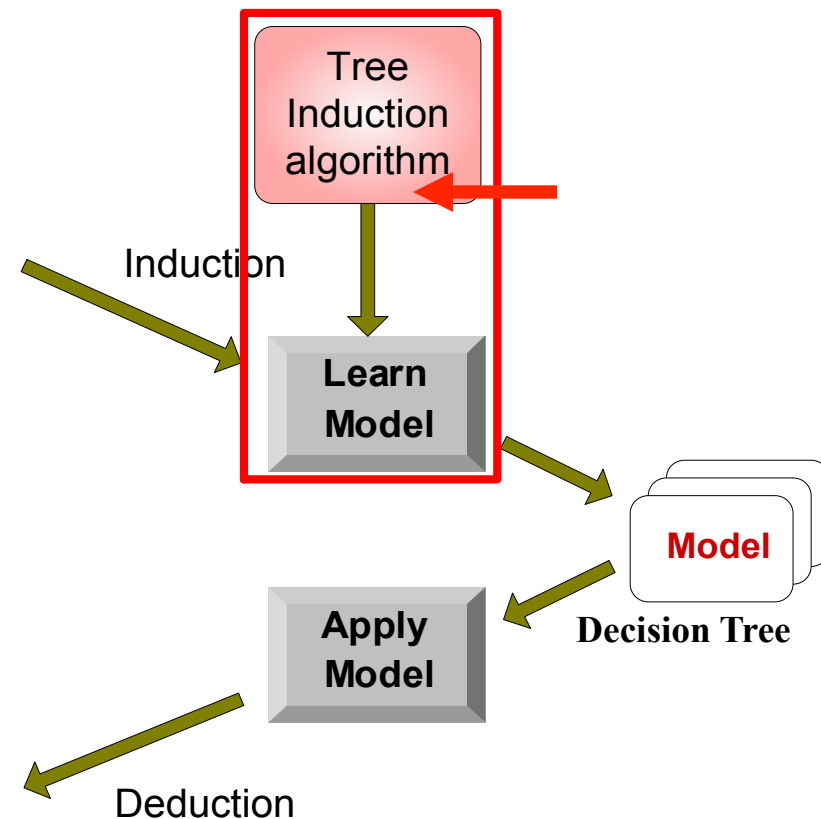
# Decision Tree Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



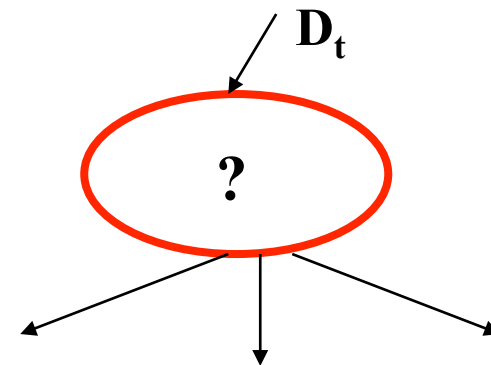
# Decision Tree Induction

- Many Algorithms:
  - Hunt's Algorithm (one of the earliest)
  - CART
  - ID3, C4.5
  - SLIQ, SPRINT

# General Structure of Hunt's Algorithm

- Let  $D_t$  be the set of training records that reach a node  $t$
- General Procedure:
  - If  $D_t$  contains records that belong the same class  $y_t$ , then  $t$  is a leaf node labeled as  $y_t$
  - If  $D_t$  is an empty set, then  $t$  is a leaf node labeled by the default class,  $y_d$
  - If  $D_t$  contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.

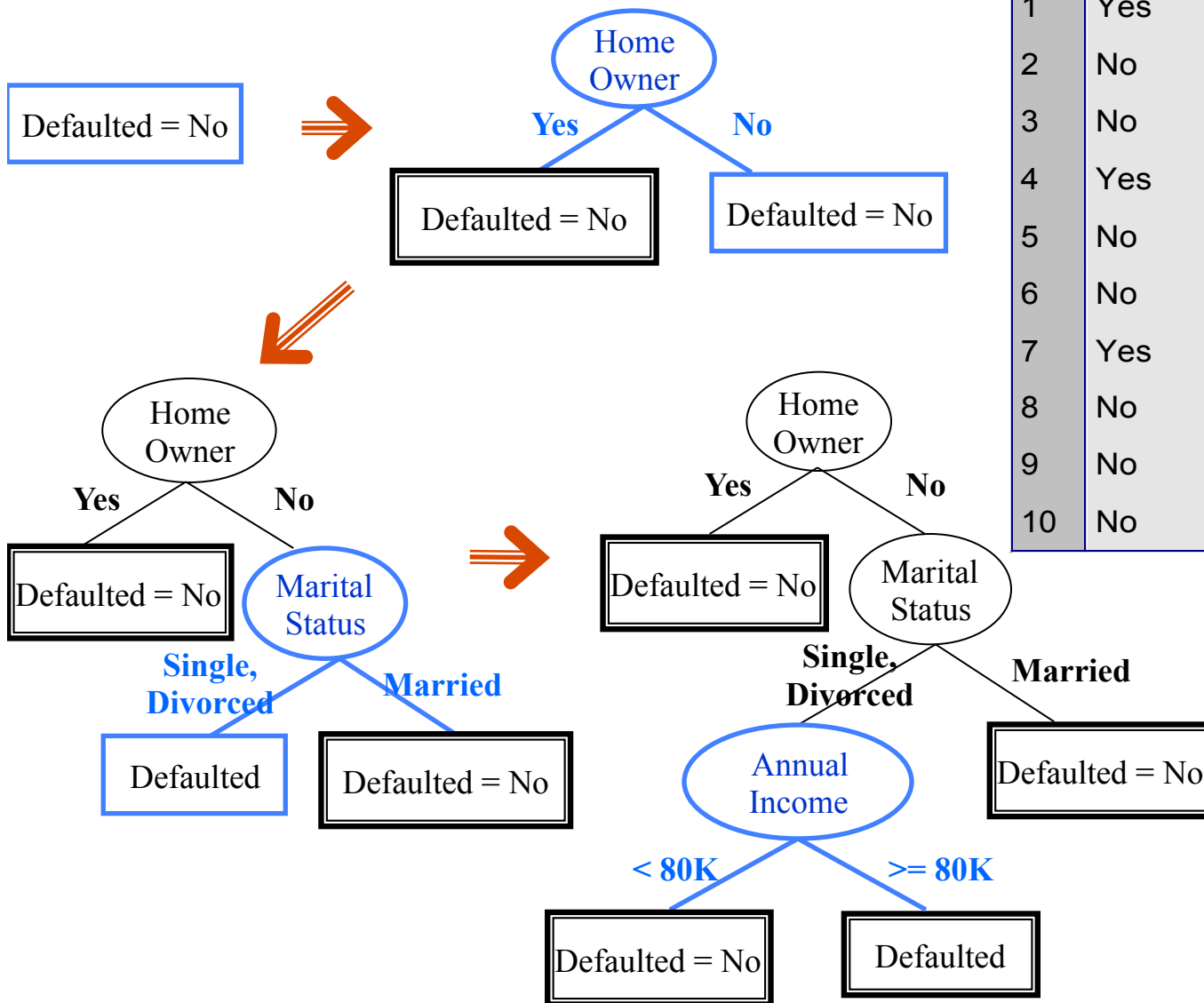
<i>Tid</i>	Home Owner	Marital Status	Annual Income	Defaulted
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes





# Hunt's Algorithm

<i>Tid</i>	Home Owner	Marital Status	Annual Income	Defaulted
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.
- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.
- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting



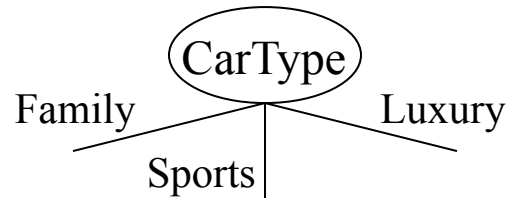
# How to Specify Test Condition?

- Depends on attribute types
  - Nominal
  - Ordinal
  - Continuous
- Depends on number of ways to split
  - 2-way split
  - Multi-way split

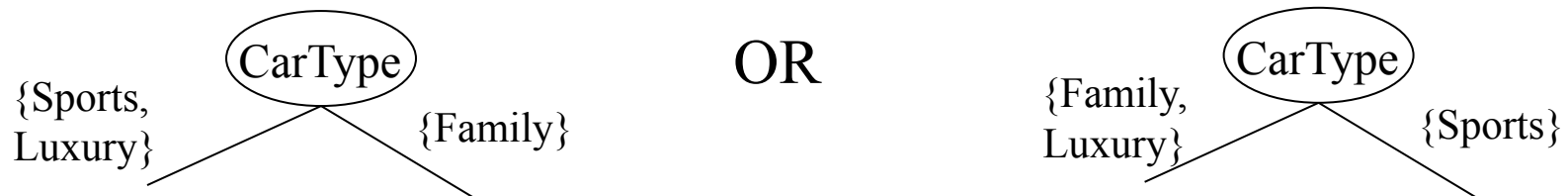


# Splitting Based on Nominal Attributes

- **Multi-way split:** Use as many partitions as distinct values.

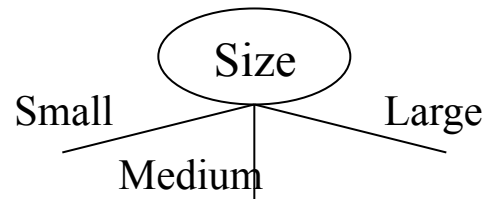


- **Binary split:** Divides values into two subsets. Need to find optimal partitioning.

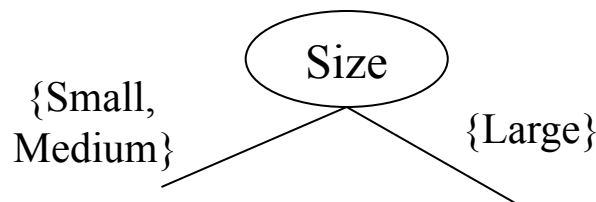


# Splitting Based on Ordinal Attributes

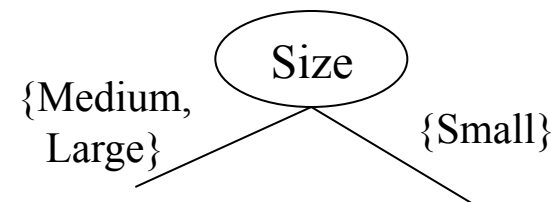
- **Multi-way split:** Use as many partitions as distinct values.



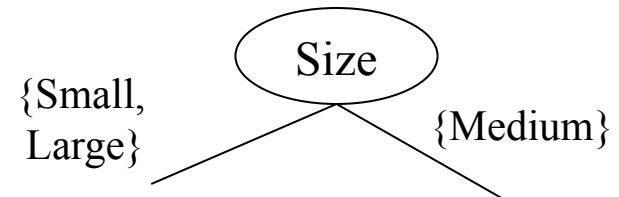
- **Binary split:** Divides values into two subsets. Need to find optimal partitioning.



OR



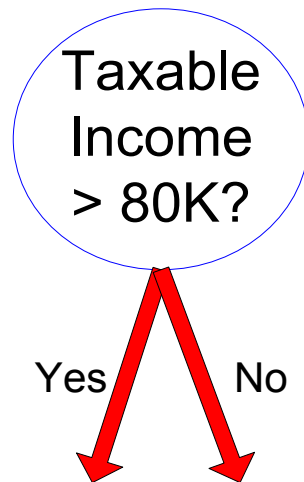
- What about this split?



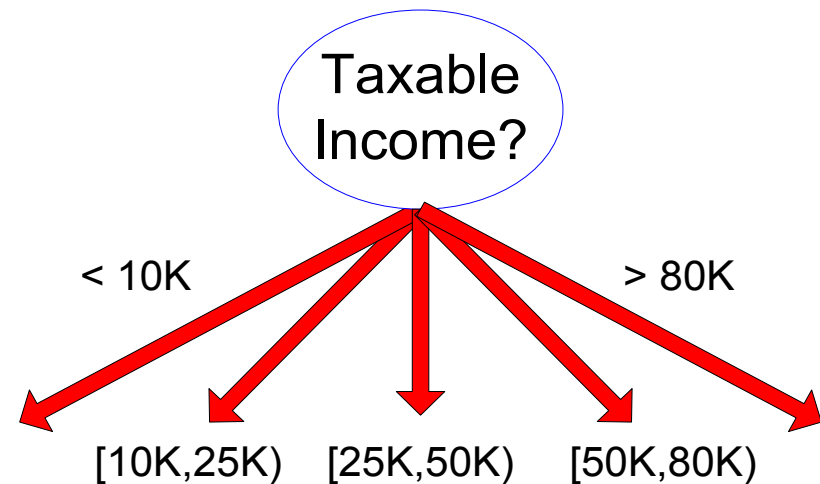
# Splitting Based on Continuous Attributes

- Different ways of handling
  - **Discretization** to form an ordinal categorical attribute
    - Static – discretize once at the beginning
    - Dynamic – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.
  - **Binary Decision**:  $(A < v)$  or  $(A \geq v)$ 
    - consider all possible splits and finds the best cut
    - can be more compute intensive

# Splitting Based on Continuous Attributes



(i) Binary split



(ii) Multi-way split

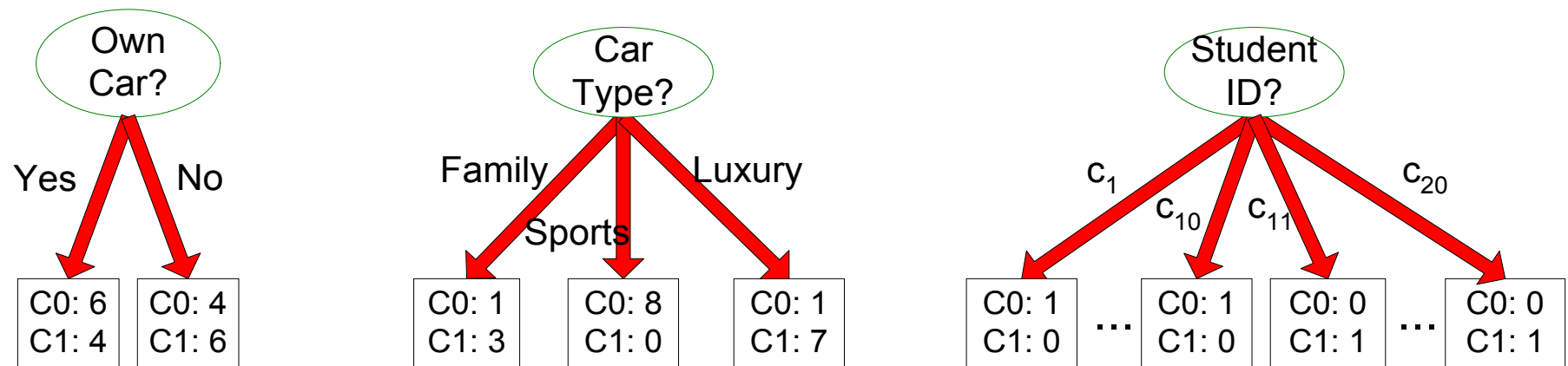


# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.
- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

# How to determine the Best Split

Before Splitting: 10 records of class 0,  
10 records of class 1



Which test condition is the best?

# How to determine the Best Split

- Greedy approach:
  - Nodes with **homogeneous** class distribution are preferred
- Need a measure of node impurity:

C0: 5
C1: 5

**Non-homogeneous,**  
**High degree of impurity**

C0: 9
C1: 1

**Homogeneous,**  
**Low degree of impurity**

---

# Measures of Node Impurity

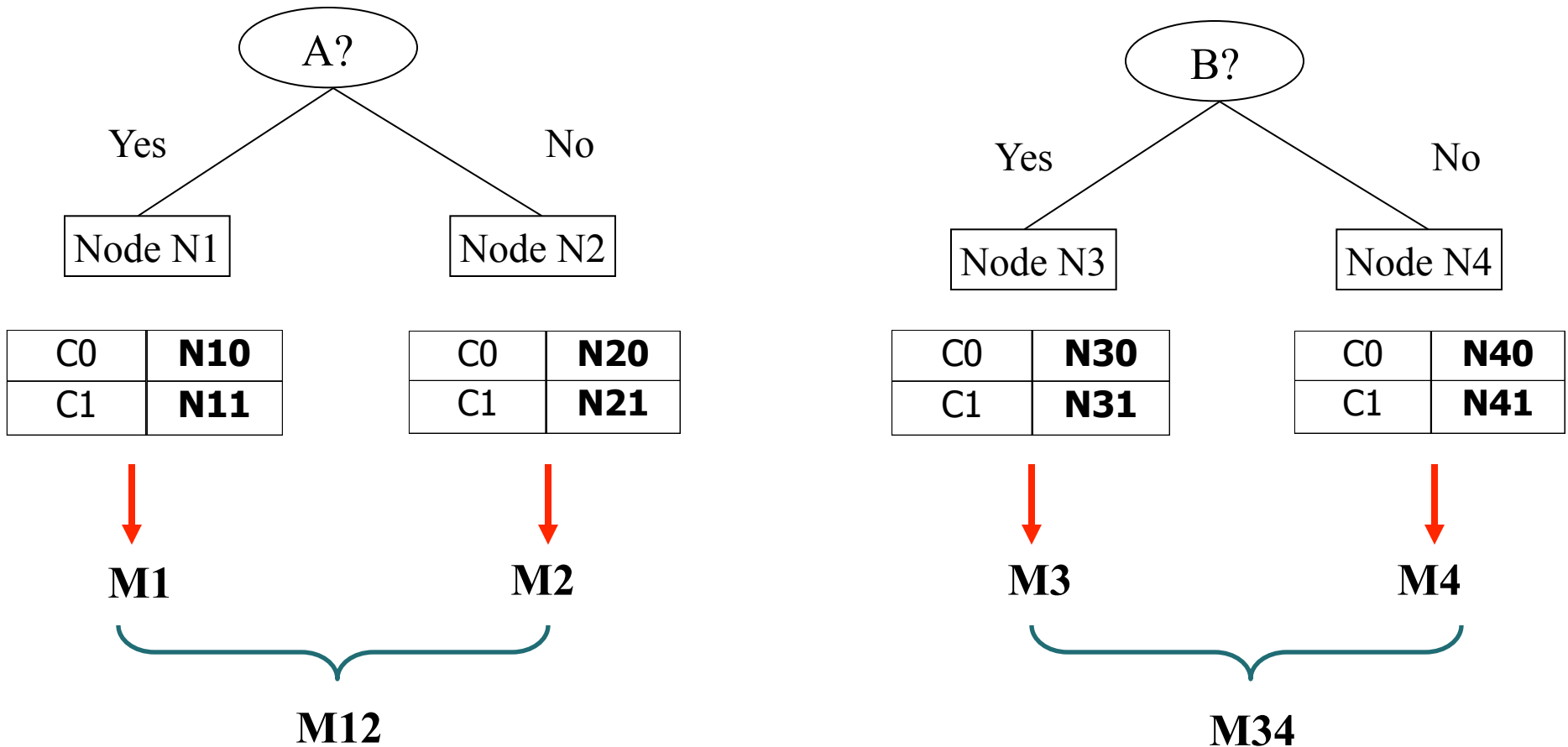
- Gini Index
- Entropy
- Misclassification error

# How to Find the Best Split

Before Splitting:

C0	<b>N00</b>
C1	<b>N01</b>

→ M0



$$\text{Gain} = M0 - M12 \text{ vs } M0 - M34$$

# Measure of Impurity: GINI

- Gini Index for a given node  $t$  :

$$GINI(t) = 1 - \sum_j [p(j | t)]^2$$

(NOTE:  $p(j | t)$  is the relative frequency of class  $j$  at node  $t$ ).

- Maximum ( $1 - 1/n_c$ ) when records are equally distributed among all classes, implying least interesting information
- Minimum (0.0) when all records belong to one class, implying most interesting information

C1	<b>0</b>
C2	<b>6</b>
<b>Gini=0.000</b>	

C1	<b>1</b>
C2	<b>5</b>
<b>Gini=0.278</b>	

C1	<b>2</b>
C2	<b>4</b>
<b>Gini=0.444</b>	

C1	<b>3</b>
C2	<b>3</b>
<b>Gini=0.500</b>	

# Examples for computing GINI

$$GINI(t) = 1 - \sum_j [p(j|t)]^2$$

C1	<b>0</b>
C2	<b>6</b>

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Gini = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$$

C1	<b>1</b>
C2	<b>5</b>

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Gini = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

C1	<b>2</b>
C2	<b>4</b>

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$Gini = 1 - (2/6)^2 - (4/6)^2 = 0.444$$

# Splitting Based on GINI

- Used in CART, SLIQ, SPRINT.
- When a node  $p$  is split into  $k$  partitions (children), the quality of split is computed as,

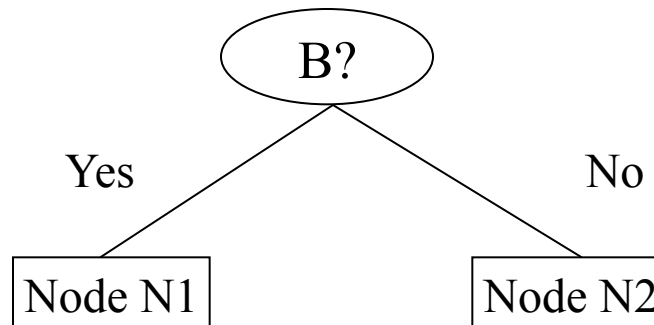
$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

where,  $n_i$  = number of records at child  $i$ ,  
 $n$  = number of records at node  $p$ .



# Binary Attributes: Computing GINI Index

- Splits into two partitions
- Effect of Weighing partitions:
  - Larger and purer partitions are sought for.



	Parent
C1	6
C2	6
<b>Gini = 0.500</b>	

$$\begin{aligned}
 \text{Gini(N1)} &= 1 - (5/7)^2 - (2/7)^2 \\
 &= 0.408
 \end{aligned}$$

$$\begin{aligned}
 \text{Gini(N2)} &= 1 - (1/5)^2 - (4/5)^2 \\
 &= 0.32
 \end{aligned}$$

	N1	N2
C1	5	1
C2	2	4
<b>Gini=0.371</b>		

$$\begin{aligned}
 \text{Gini(Children)} &= 7/12 * 0.408 + \\
 &\quad 5/12 * 0.32 \\
 &= 0.371
 \end{aligned}$$