



香港城市大學
City University of Hong Kong

專業 創新 胸懷全球
Professional • Creative
For The World

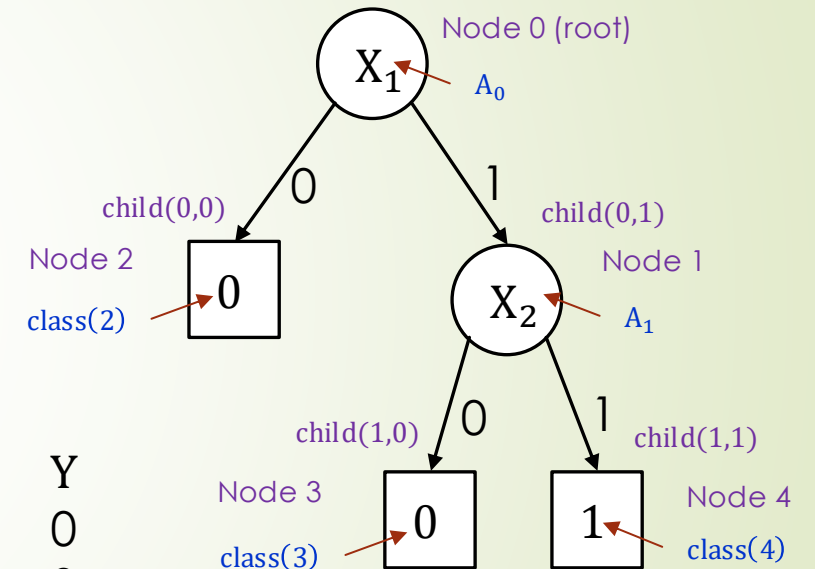
Classification: Decision Tree Induction

C5483 Data Warehousing and Data Mining

What is a decision tree?

- Internal nodes t (circle)
 - Label A_t (splitting criterion)
 - For each $A_t = j$ (outcome), an edge to $\text{child}(t, j)$ (child node)
- Leaf nodes (square)
 - label $\text{class}(t)$ (decision)

	X_1	X_2	Y
1.	0	0	0
2.	0	1	0
3.	1	0	0
4.	1	1	1
	x_1	x_2	y



3

How to classify?

Trace from root to leaves

Input: feature vector x

Output: predicted class \hat{y}

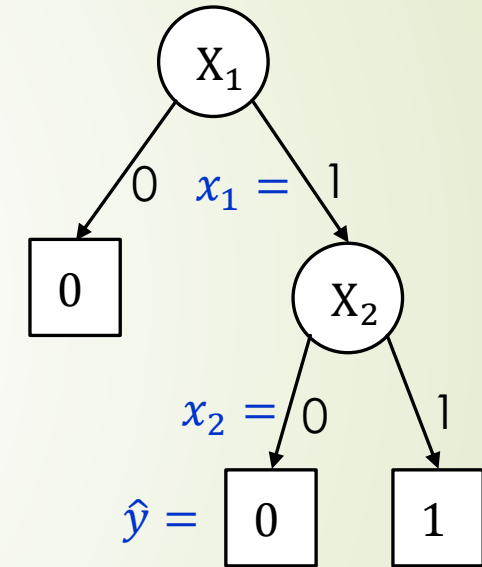
$t \leftarrow \text{root}$

while t is not a leaf

$t \leftarrow \text{child}(t, j)$ where $A_t = j$ for x

$\hat{y} \leftarrow \text{class}(t)$

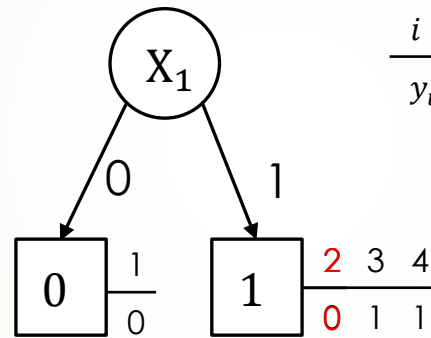
$$x = (x_1, x_2) = (1, 0)$$



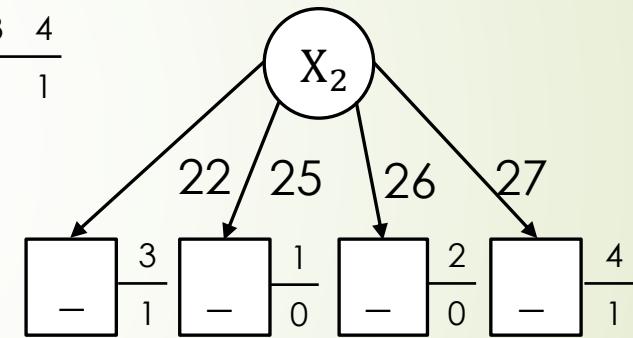
How to build a decision stump?

- A decision stump is a decision tree with depth ≤ 1 .

	X_1	X_2	Y
1.	0	25	0
2.	1	26	0
3.	1	22	1
4.	1	27	1



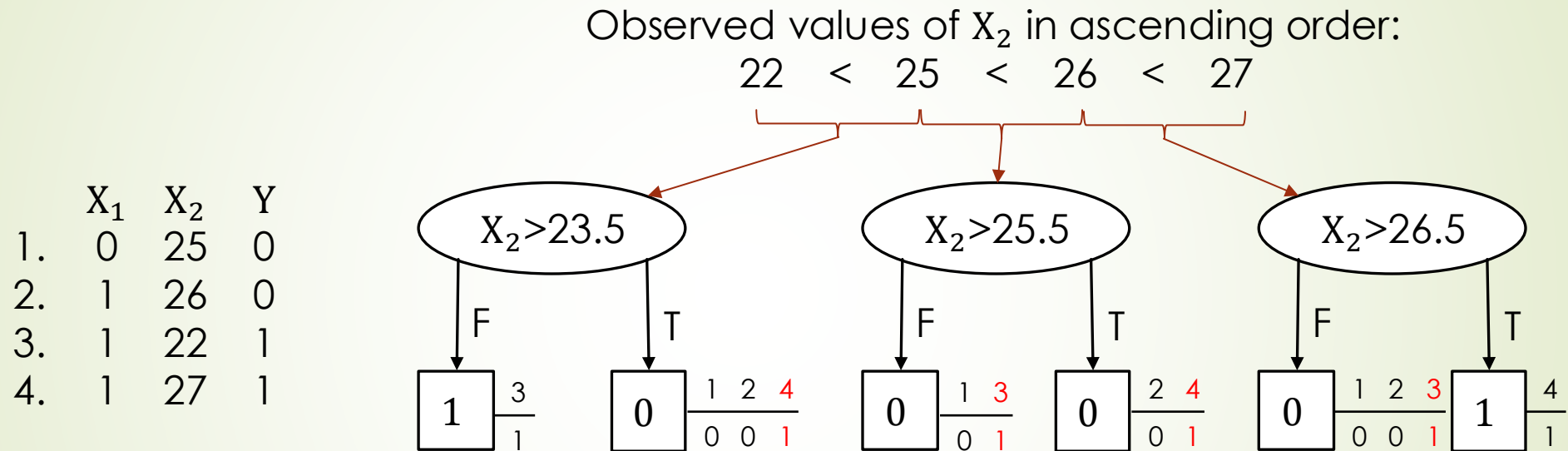
error rate = _____



error rate = _____

- Choose a splitting attribute.
- Use majority voting to determine $\text{class}(t)$.
- Which decision stump is better? Left/right because of _____.

Binary splits for numeric attributes



- **c**_____ **m**____-points as **s**_____ points.
- Which is/are the best split(s)? _____
- How to build a tree instead of a stump? R_____ly split (d_____ and c_____).

How to build a decision tree?

Greedy algorithm (See [Han11 Fig 8.3](#) for the full version)

Input: training data D

Output: root node of the decision tree

function Split(D)

 create node t

$A_t \leftarrow$ a good criterion A to split D

 if A_t is not null

 for each outcome j of A_t

$D_j \leftarrow \{(x, y) \in D \mid x \text{ satisfies } A_t = j\}$

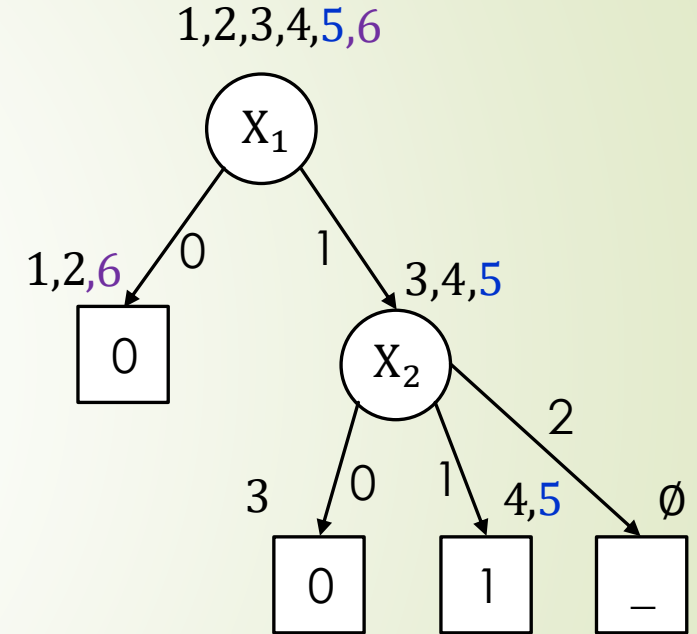
 child(t, j) \leftarrow _____ if $D_j \neq \emptyset$
 else _____

 else # further splitting is not useful when _____

 class(t) \leftarrow _____

 return t

	X_1	X_2	Y
1.	0	0	0
2.	0	1	0
3.	1	0	0
4.	1	1	1
5.	1	1	0
6.	0	2	0



(a) $\arg \max_k |\{(x, y) \in D \mid y = k\}|$

(b) new node t' with class(t') $\leftarrow \arg \max_k |\{(x, y) \in D \mid y = k\}|$

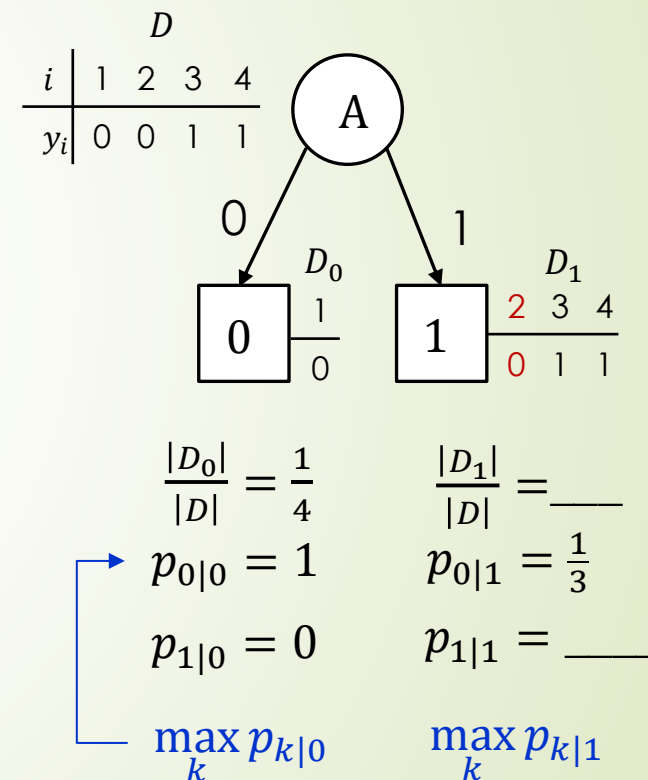
(c) Split(D_j)

How to find good splitting attribute?

- Given the data D to split, choose the splitting attribute A that minimizes e _____ of decision stump by A .
- What is the precise formula?

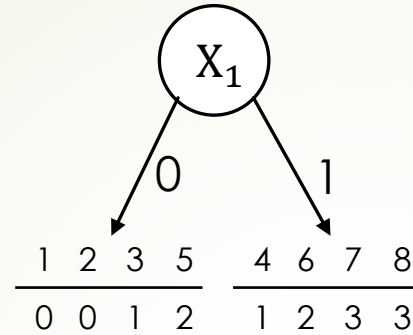
$$\begin{aligned} \text{Misclass}_A(D) &:= \sum_j \frac{|D_j|}{|D|} \left(1 - \max_k p_{k|j} \right) \\ &= 1 - \underbrace{\sum_j \frac{|D_j|}{|D|} \max_k p_{k|j}}_{\text{accuracy}} \end{aligned}$$

- D_j : set $\{(x, y) \in D \mid A = j \text{ for } x\}$ of tuples in D satisfying $A = j$.
- $p_{k|j}$: fraction $\frac{|\{(x, y) \in D_j \mid y = k\}|}{|D_j|}$ of tuples in D_j belonging to class k .

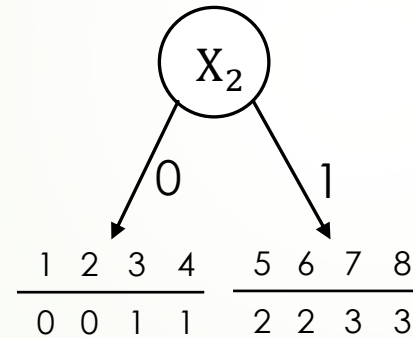


$\text{Misclass}_A(D) =$ _____

	X_1	X_2	X_3	Y
1.	0	0	0	0
2.	0	0	0	0
3.	0	0	1	1
4.	1	0	1	1
5.	0	1	0	2
6.	1	1	0	2
7.	1	1	1	3
8.	1	1	1	3



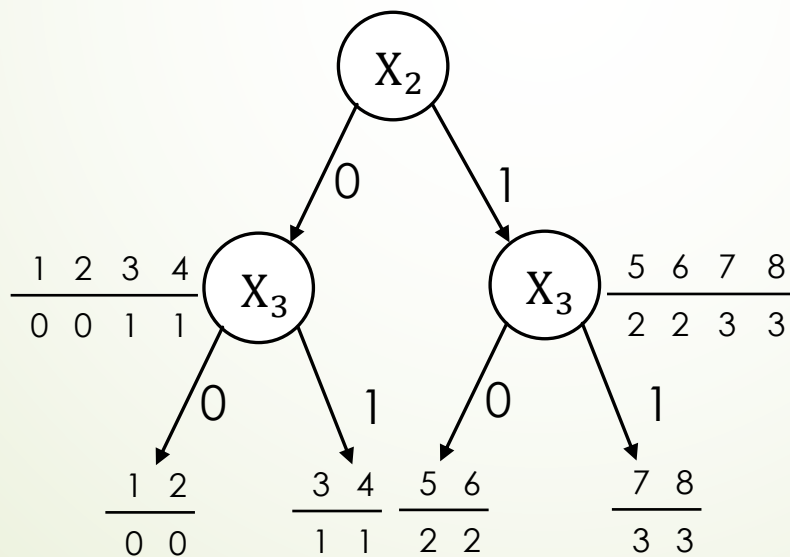
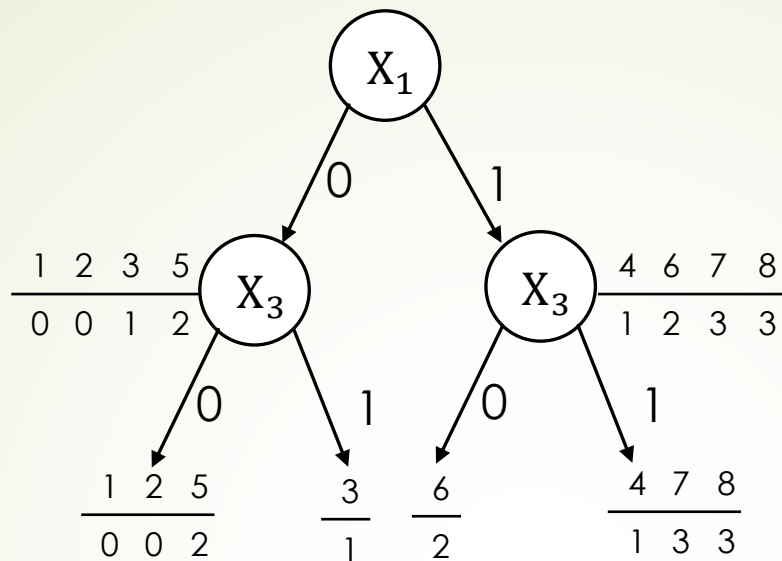
$$\text{Misclass}_{X_1}(D) = \underline{\hspace{2cm}}$$



$$\text{Misclass}_{X_2}(D) = \underline{\hspace{2cm}}$$

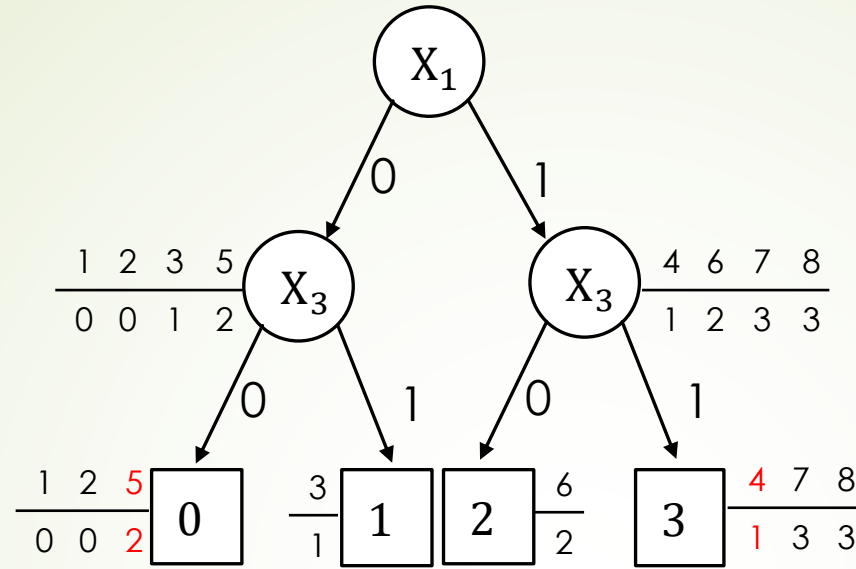
- What is the best splitting attribute? X_1 / X_2 / same

	X_1	X_2	X_3	Y
1.	0	0	0	0
2.	0	0	0	0
3.	0	0	1	1
4.	1	0	1	1
5.	0	1	0	2
6.	1	1	0	2
7.	1	1	1	3
8.	1	1	1	3

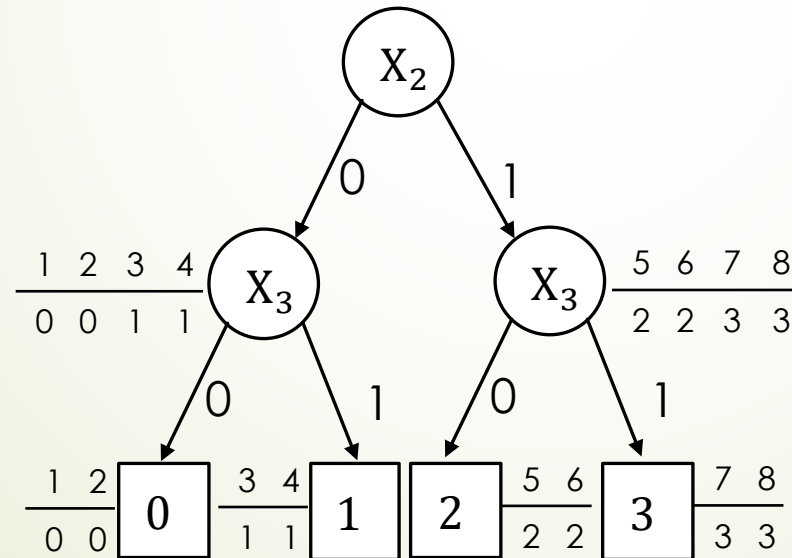


Further split on X_3

	X_1	X_2	X_3	Y
1.	0	0	0	0
2.	0	0	0	0
3.	0	0	1	1
4.	1	0	1	1
5.	0	1	0	2
6.	1	1	0	2
7.	1	1	1	3
8.	1	1	1	3



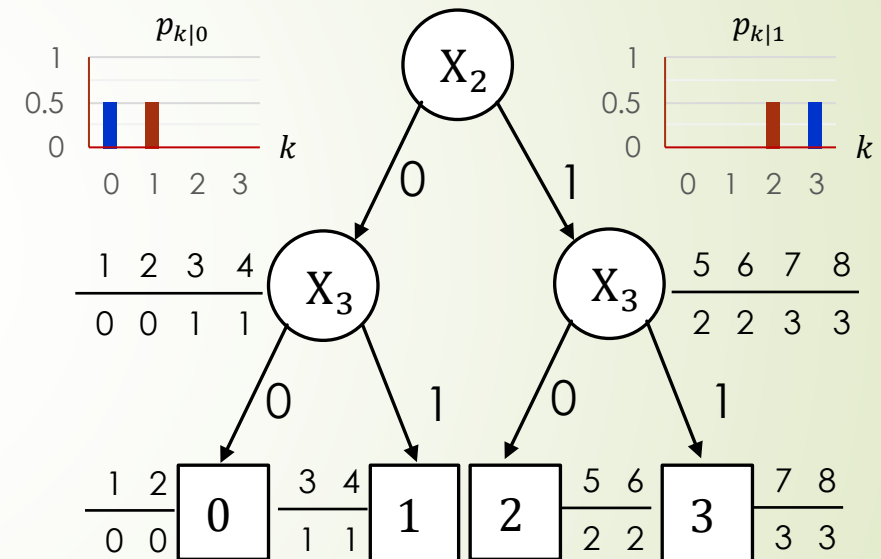
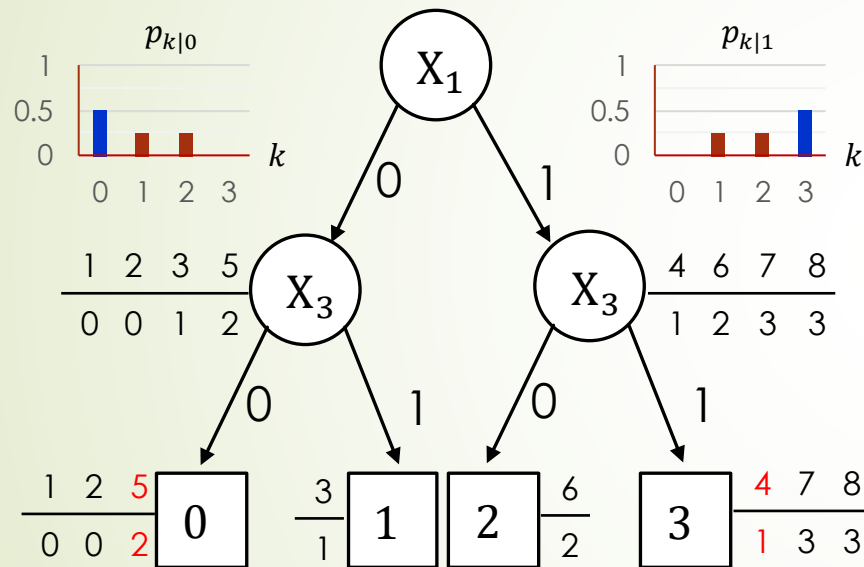
Error rate = _____



Error rate = _____

Issue of greedy algorithm

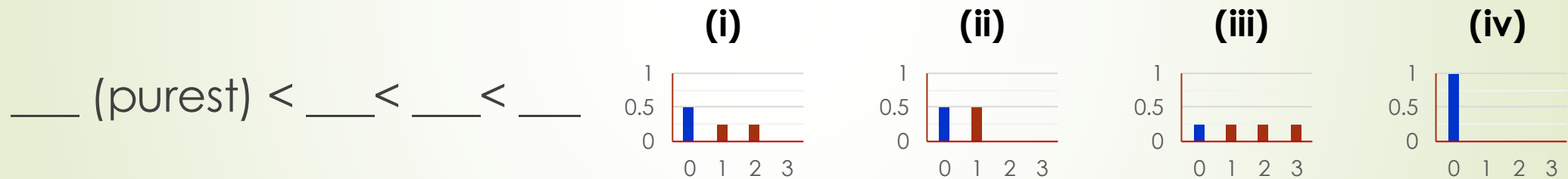
- **Locally optimal** split may not be **g**_____ **optimal**.



- Why splitting on X_1 is not good? Child nodes of X_1 are **less p**_____.
- Why misclassification rate fails?
It neglects the **distribution of the class values of m**_____ **instances**.

How to remain greedy but not myopic

- Find better **i** _____ **measures** than misclassification rate.
- How to measure impurity?
- E.g., order the following distributions in ascending order of impurities:



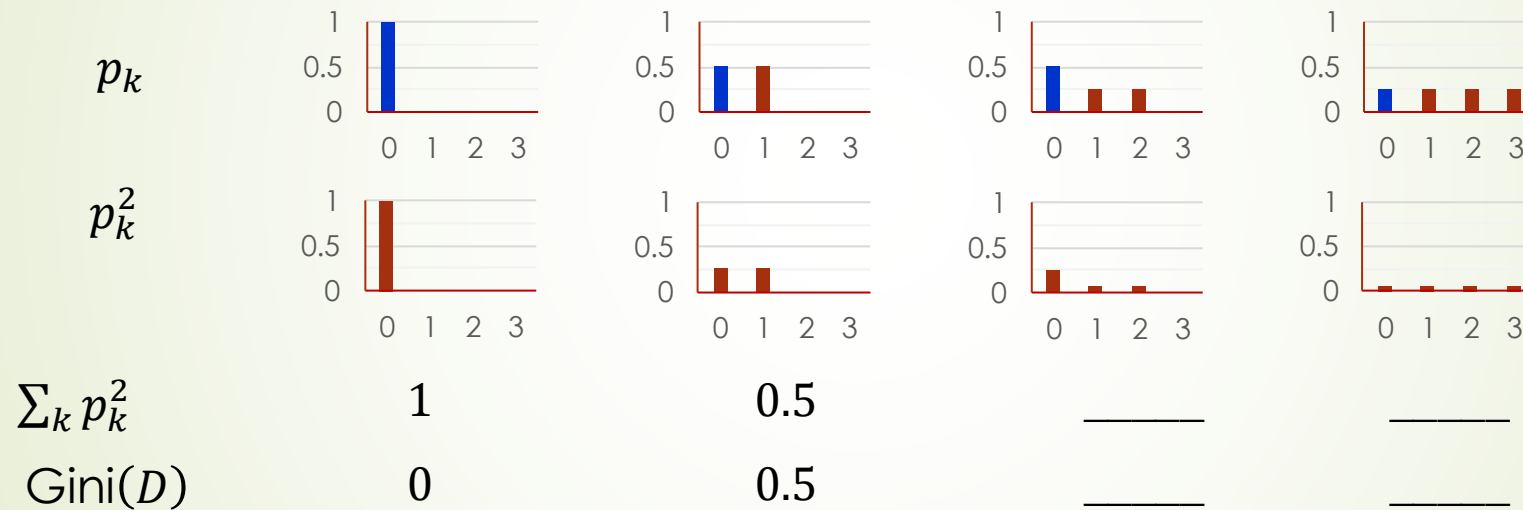
- Given a distribution p_k of the class values of D , how to define a non-negative function of p_k 's that respect the above ordering?

$1 - \max_k p_k$ works? Yes/No

$1 - \sum_k p_k$ works? Yes/No

Gini impurity index

$$\text{Gini}(D) := g(p_0, p_1, \dots) := \sum_k p_k(1 - p_k) = 1 - \sum_k p_k^2$$



Why it works?

➤ $g(p_0, p_1, \dots) \geq 0$. Equality iff $\forall k, p_k \in \{0,1\}$. Why?

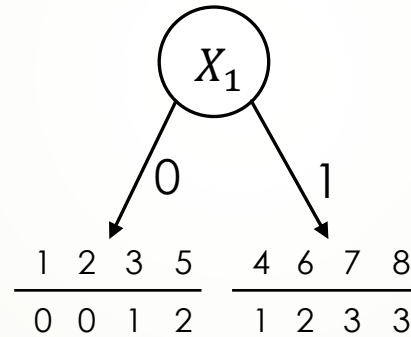
➤ $g(p_0, p_1, \dots, p_n) \leq 1 - \frac{1}{n}$. Equality iff $p_k = \frac{1}{n}$. Why?

Finding the best split using Gini impurity

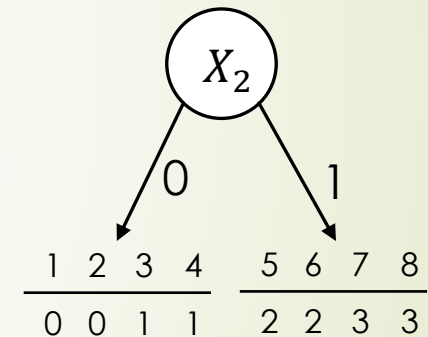
- Minimize the Gini impurity given A:

$$\text{Gini}_A(D) := \sum_j \frac{|D_j|}{|D|} \text{Gini}(D_j)$$

	X_1	X_2	X_3	Y
1.	0	0	0	0
2.	0	0	0	0
3.	0	0	1	1
4.	1	0	1	1
5.	0	1	0	2
6.	1	1	0	2
7.	1	1	1	3
8.	1	1	1	3



$$\begin{aligned} \text{Gini}(D_0) = \text{Gini}(D_1) &= 0.625 \\ &= \text{Gini}_A(D) \end{aligned}$$



$$\begin{aligned} \text{Gini}(D_0) = \text{Gini}(D_1) &= \text{_____} \\ &= \text{Gini}_A(D) \end{aligned}$$


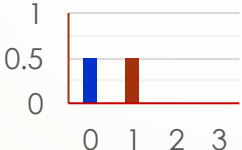
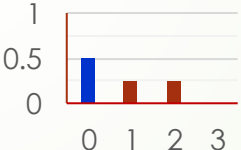
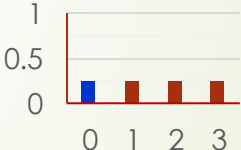
- What is the best splitting attribute? X_1 / X_2 / same

An impurity measure from information theory

Shannon's entropy

$$\begin{aligned}\text{Info}(D) &:= h(p_1, p_2, \dots) \\ &:= \sum_k p_k \log \frac{1}{p_k} = - \sum_{k:p_k > 0} p_k \log p_k\end{aligned}$$

- Measured in bits with base-2 logarithm. Why?
- $0 \log 0$ is regarded as $\lim_{p \rightarrow 0} p \log p$ even though $\log 0$ is undefined.

p_k					
$p_k \log_2 \frac{1}{p_k}$	0,0,0,0	$a, a, 0, 0$	$a, b, b, 0$	b, b, b, b	$a = \underline{\hspace{1cm}}$ $b = \underline{\hspace{1cm}}$
$\text{Info}(D)$	0	$\underline{\hspace{1cm}}$	$\underline{\hspace{1cm}}$	$\underline{\hspace{1cm}}$	

Why it works?

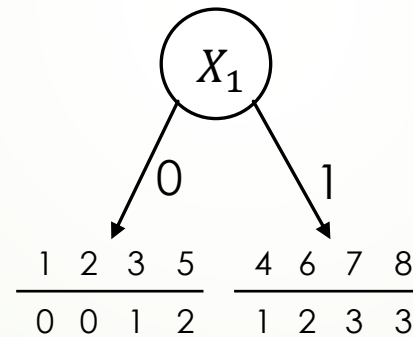
- $h(p_0, p_1, \dots) \geq 0$. Equality iff $\forall k, p_k \in \{0,1\}$. Why?
- $h(p_0, p_1, \dots, p_n) \leq \log_2 n$. Equality iff $p_k = \underline{\hspace{2cm}}$. Why?

Finding the best split by conditional entropy

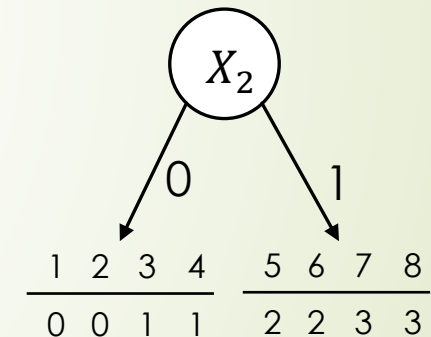
- Minimize the entropy given A,

$$\text{Info}_A(D) := \sum_j \frac{|D_j|}{|D|} \text{Info}(D_j)$$

	X_1	X_2	X_3	Y
1.	0	0	0	0
2.	0	0	0	0
3.	0	0	1	1
4.	1	0	1	1
5.	0	1	0	2
6.	1	1	0	2
7.	1	1	1	3
8.	1	1	1	3



$$\begin{aligned} \text{Info}(D_0) = \text{Info}(D_1) &= 1.5 \\ &= \text{Info}_A(D) \end{aligned}$$



$$\begin{aligned} \text{Info}(D_0) = \text{Info}(D_1) &= \text{_____} \\ &= \text{Info}_A(D) \end{aligned}$$

- What is the best splitting attribute? X_1 / X_2 / same

Which impurity measure is used?

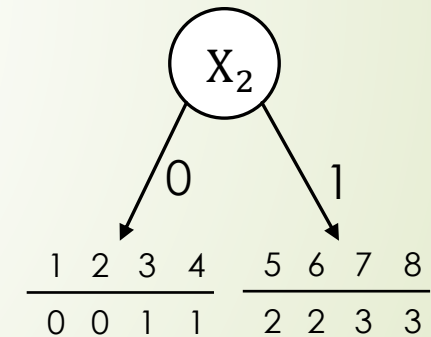
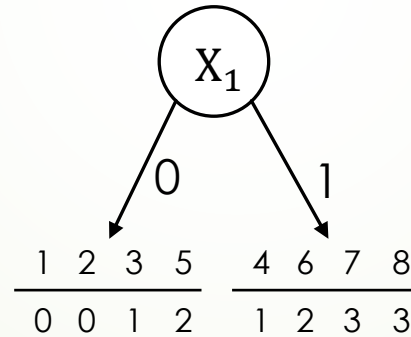
- **ID3** (Iterative Dichotomiser 3) maximizes

$\text{Gain}_A(D) := \text{Info}(D) - \text{Info}_A(D)$ (information gain or mutual information)

- **CART** (Classification and Regression Tree)

$\Delta\text{Gini}_A(D) := \text{Gini}(D) - \text{Gini}_A(D)$ (Drop in Gini impurity)

	X_1	X_2	X_3	Y
1.	0	0	0	0
2.	0	0	0	0
3.	0	0	1	1
4.	1	0	1	1
5.	0	1	0	2
6.	1	1	0	2
7.	1	1	1	3
8.	1	1	1	3



$\text{Info}(D) = 2$
 $\text{Gini}(D) = 0.75$

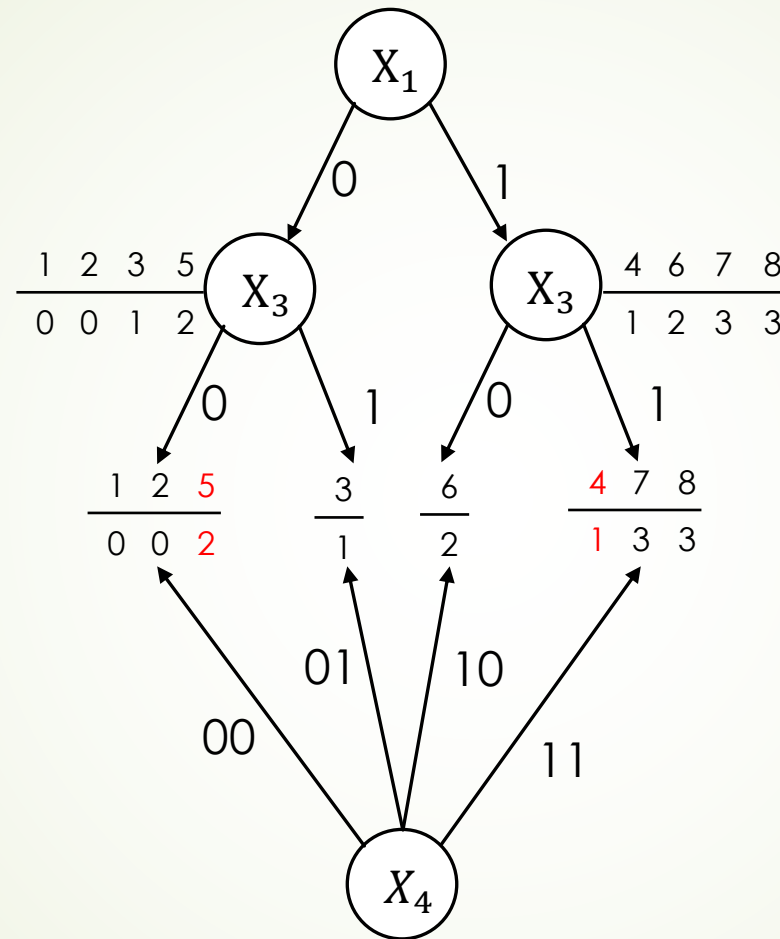
$\text{Info}_{X_1}(D) = 1.5$ $\text{Gain}_{X_1}(D) = 0.5$
 $\text{Gini}_{X_1}(D) = 0.625$ $\Delta\text{Gini}_{X_1}(D) = 0.125$

$\text{Info}_{X_2}(D) = 1$ $\text{Gain}_{X_2}(D) = \underline{\hspace{2cm}}$
 $\text{Gini}_{X_2}(D) = 0.5$ $\Delta\text{Gini}_{X_2}(D) = \underline{\hspace{2cm}}$

- What is the best splitting attribute? X_1 / X_2 / same

	X_1	X_2	X_3	X_4	Y
1.	0	0	0	00	0
2.	0	0	0	00	0
3.	0	0	1	01	1
4.	1	0	1	11	1
5.	0	1	0	00	2
6.	1	1	0	10	2
7.	1	1	1	11	3
8.	1	1	1	11	3

$$X_4 := X_1 \circ X_3$$



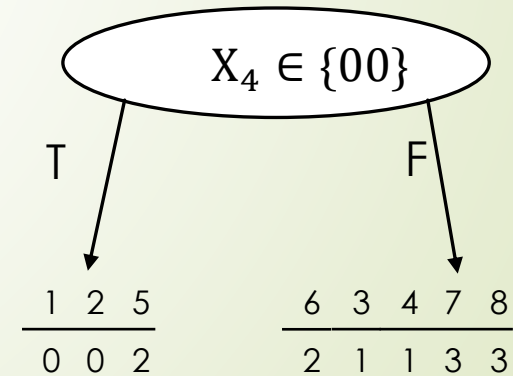
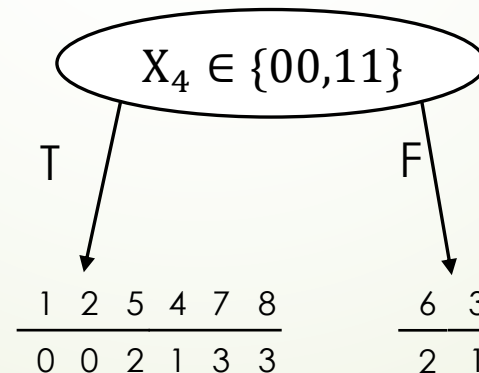
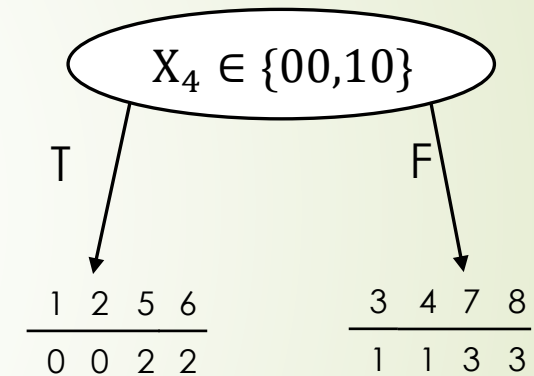
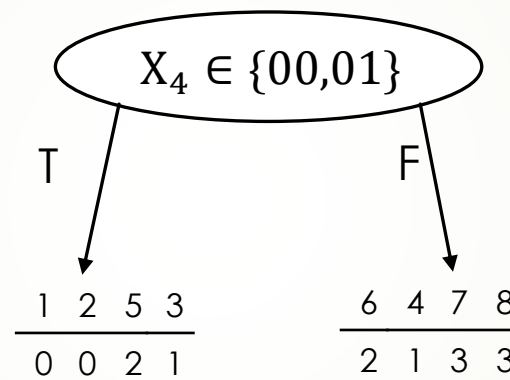
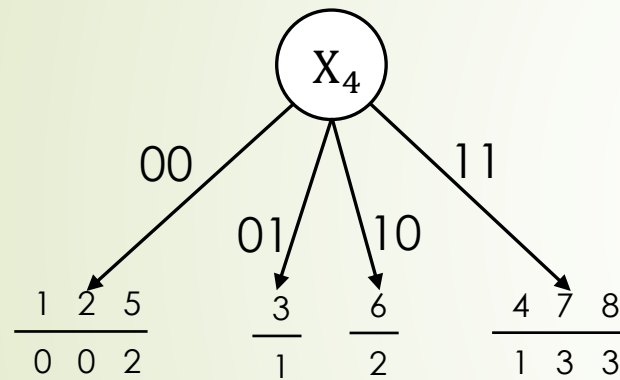
➤ Is X_4 a good splitting attribute? Yes/No.

Bias towards attributes with many outcomes

- An attribute with more outcomes tends to
 - reduce impurity more but
 - result in more comparisons.
- Issues: Such attribute may not *minimize impurity per comparison*.
- Remedies?

Binary split also for nominal attributes

- CART uses a **s**_____ S to generate a **binary split** (whether $A \in S$).
- The number of outcomes is therefore limited to _____.



...

$\max_S \Delta \text{Gini}_{X_4 \in S}(D) = \underline{\hspace{2cm}}$
 achieved by $S = \underline{\hspace{2cm}}$

Normalization by split information

- C4.5/J48 allows **m** _____ split but uses **information gain ratio**

$$\frac{\text{Gain}_A(D)}{\text{SplitInfo}_A(D)} \quad \text{where} \quad \text{SplitInfo}_A(D) = \sum_j \frac{|D_j|}{|D|} \log_2 \frac{1}{|D_j|/|D|}.$$

- $\text{SplitInfo}_A(D)$ is the entropy of _____ because _____.
- Attributes with many outcomes tend to have smaller/larger $\text{SplitInfo}_A(D)$.

How to avoid overfitting

- **P__-pruning**: Limit the size of the tree as we build it. E.g.,
 - Ensure each node is supported by enough examples. (C4.5: minimum number of objects.)
 - Split only if we are confident enough about the improvement. (C4.5: confidence factor.)
- **P__-pruning**: Reduce the size of the tree after we build it. E.g.,
 - Contract leaf nodes if complexity outweighs the risk. (CART: **cost-complexity pruning**)

References

- 8.1 Basic Concepts
- 8.2 Decision Tree Induction
- Optional readings
 - https://en.wikipedia.org/wiki/C4.5_algorithm
 - [Cover, T., & Thomas, J. \(2006\). *Elements of information theory* \(2nd ed.\). Hoboken, N.J.: Wiley-Interscience.](#) Chapter 1 and 2.