

Machine Learning: Classification Techniques

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Prajwol Sangat



Last week

Data Transformer, Estimators, Pipelines
Feature Selection and Extraction

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

This week

Classification Algorithms

Decision Tree

Random Forest

Assignment Project Exam Help

DEMO

<https://powcoder.com>

Add WeChat powcoder

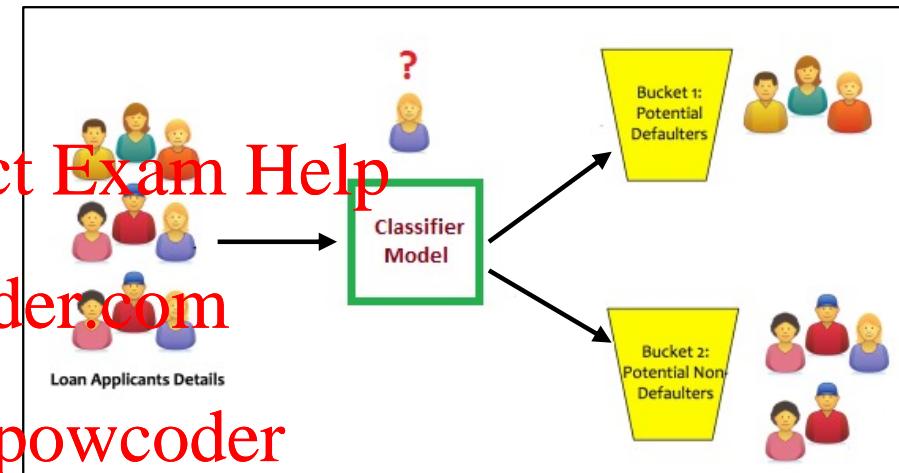
Classification

Predictive Data Modeling

A classifier model needs to be created using training dataset

After the classifier is created, classification is the process of assigning new instances from the testing dataset to predefined classes

The label for each class is predefined



Classification

Classifiers can be:

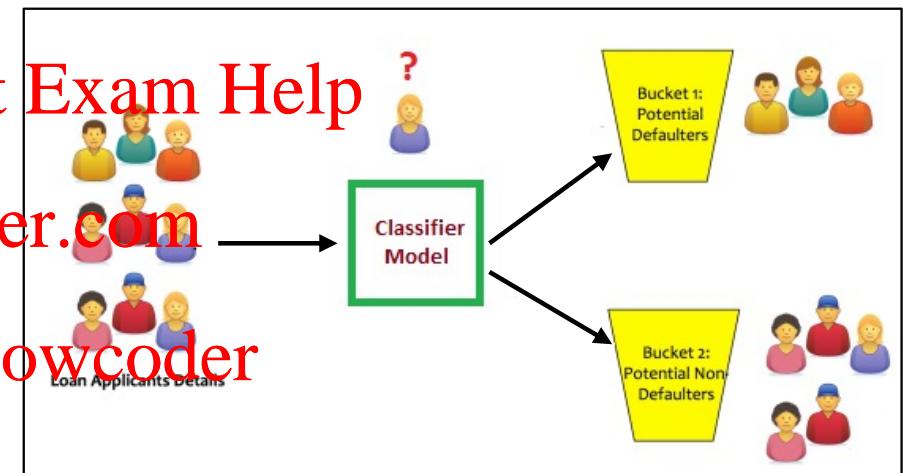
- *Binary classifier*
- *Multi-Class classifiers*

Binary classifiers: Classification with only 2 distinct classes or with 2 possible outcomes

Example: classification of spam email and non-spam email, potential defaulter and non defaulter

Multi-Class classifiers: Classification with more than two distinct classes

Example: classification of types of animals, classification of books into categories.



Classification Algorithms

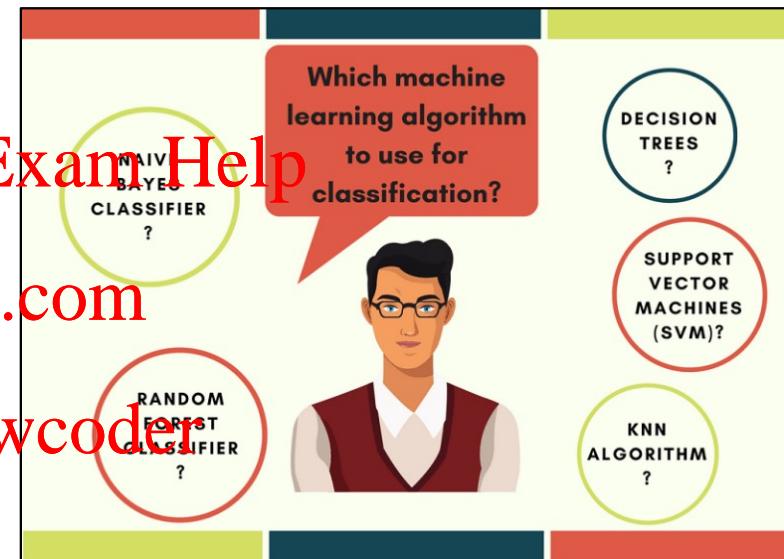
There are several types of classification algorithms in Machine Learning:

- Decision Trees
- Random Forest
- Logistic Regression
- And Many More.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



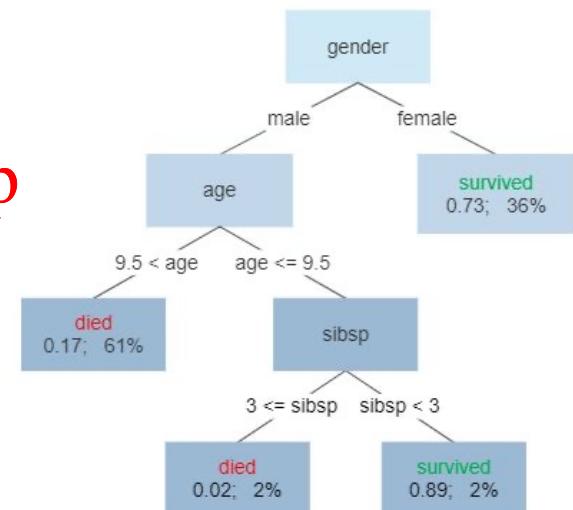
Decision Trees

- A tree-like predictive model for decision making
- In DTs, a record/sample which falls into a certain class or category is identifiable through its features/attributes.
- It splits samples into two or more homogeneous sets (leaves) based on the most significant attributes (predictors)

Assignment Project Exam Help
<https://powcoder.com>

Samples	Features/Attributes			Class
	gender	age	sibsp	
Person 1	male	30	1	died
Person 2	female	20	2	survived

Survival of passengers on the Titanic



Example: Titanic dataset

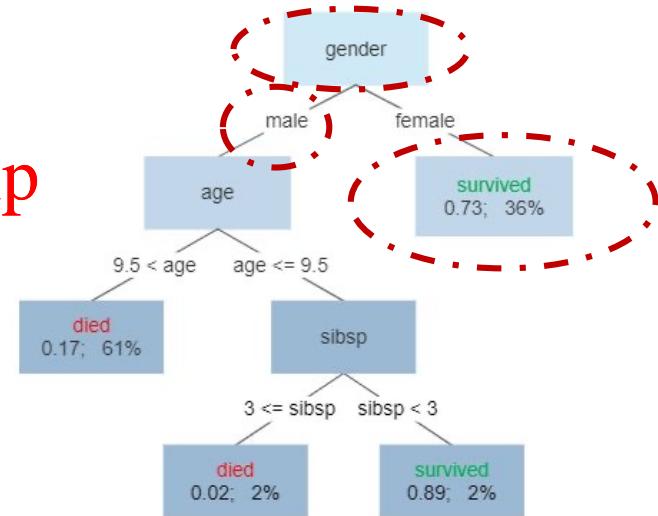
Decision Trees

- Each **internal node** represents a "test" on an attribute (e.g. gender)
- Each **branch** corresponds to attribute values (outcome of test) – e.g. male or female
- Each **leaf/terminal node** assigns class label (e.g., died or survived)

Assignment Project Exam Help

<https://powcoder.com>
Add WeChat powcoder

Survival of passengers on the Titanic



Example: Titanic dataset

Decision Tree Algorithm

Common terms used with Decision trees:

Root Node: It represents entire population or sample and this further gets divided into two or more homogeneous sets.

Splitting: It is a process of dividing a node into two or more sub-nodes.

Decision Node: When a sub-node splits into further sub-nodes, then it is called decision node.

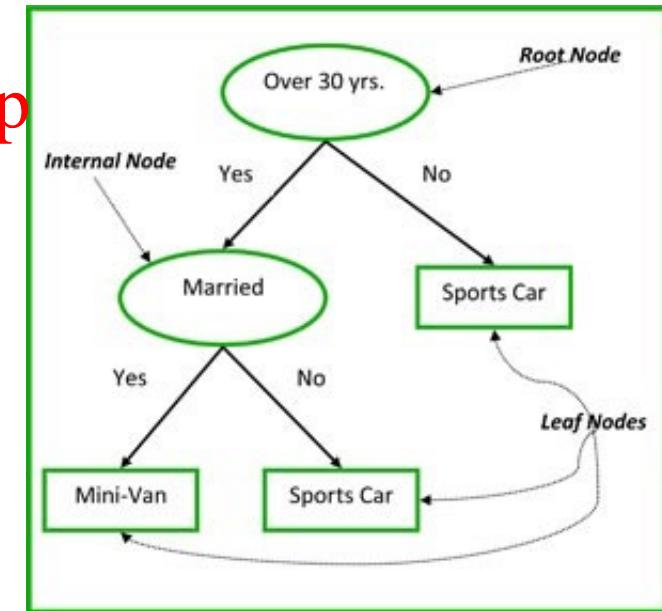
Assignment Project Exam Help
<https://powcoder.com>

Leaf/ Terminal Node: Nodes do not split is called Leaf or Terminal node.

Pruning: When sub-nodes of a decision node is removed, this process is called pruning (an opposite process of splitting).

Branch / Sub-Tree: A sub section of entire tree is called branch or sub-tree.

Parent and Child Node: A node, which is divided into sub-nodes is called parent node of sub-nodes whereas sub-nodes are the child of parent node.



Decision Tree Algorithm

Supervised Learning – need output labels to build a DT

Constructing a DT is generally a recursive process

- Initialization: All training data at the root node
- Partition training data recursively by choosing one attribute at a time
- Repeat process for partitioned dataset
- Stopping criteria: When all training data in each partition have same target class

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The most common approach in building a decision tree:

ID3 (Iterative Dichotomiser 3) → uses **Entropy function** and **Information gain** as metrics to construct a DT.

ID3 (Iterative Dichotomiser 3)

ID3 (Iterative Dichotomiser 3) was developed in 1986 by Ross Quinlan.

Assignment Project Exam Help

The algorithm creates **a multiway tree**, finding for each node (i.e. in a greedy manner) [the categorical feature that will yield the largest information gain for categorical targets.](https://powcoder.com)

Add WeChat powcoder

Trees are grown to their maximum size and then a pruning step is usually applied to improve the ability of the tree to generalise to unseen data.

Entropy

- Measure of uncertainty or randomness in data
- Informs the predictability of an event
 - Low value -> Less uncertainty, high value -> high uncertainty

Less homogeneous

Play Basketball	
Yes	No
9	5

More homogeneous

Play Basketball	
Yes	No
13	1

Assignment Project Exam Help

$$H(S) = \sum_{i=1}^n p_i \log \frac{1}{p_i}$$

p_i - Probability of event i
 n - Number of events

<https://powcoder.com>

$$H(\text{Play_basketball}) = p(\text{yes}) \log \frac{1}{p(\text{yes})} + p(\text{no}) \log \frac{1}{p(\text{no})}$$
$$= -\left(\frac{9}{14} \log \frac{9}{14}\right) - \left(\frac{5}{14} \log \frac{5}{14}\right)$$
$$= 0.2831$$

$$H(\text{Play_basketball}) = -\left(\frac{13}{14} \log \frac{13}{14}\right) - \left(\frac{1}{14} \log \frac{1}{14}\right)$$
$$= 0.1115$$

If samples are completely homogeneous, the entropy is zero

Information gain

- IG for a set S is change in entropy after deciding on a attribute A.
- It computes difference between entropy before split and average entropy after split of the dataset based on an attribute A
- Used to decide which attributes are more relevant in ID3 algorithm

$$IG(S, A) = H(S) - H(S, A)$$

$$= H(S) - \sum_{i \in \text{values}(A)} p_i H(S_i)$$

<https://powcoder.com>

Entropy before
(on entire set A) Entropy after a decision
based on A

Add WeChat powcoder

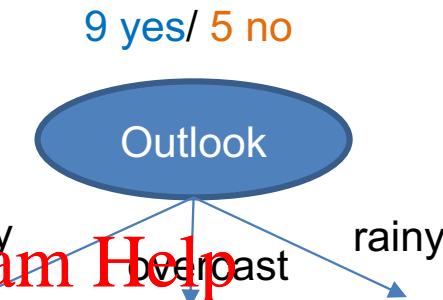
S_i Subset/partition of data after splitting S

Example

		Play Basketball		
		Yes	No	
Outlook	Sunny	3	2	5
	Overcast	4	0	4
	Rainy	2	3	5
				14

Play Basketball	
Yes	No
9	5

$$H(S) = p(yes) \log \frac{1}{p(yes)} + p(no) \log \frac{1}{p(no)} \\ = 0.2831$$



Assignment Project Exam Help

<https://powcoder.com>

$$H(S,A) = p(sunny)H(S_{sunny}) + p(overcast)H(S_{overcast}) + p(rain)H(S_{rain})$$

$$\begin{aligned} &= \frac{5}{14}(0.2922) + \frac{4}{14}(0) + \frac{5}{14}(0.2922) \\ &= 0.2087 \end{aligned}$$

In ID3 algorithm, we select attribute with the highest gain to be the node in the tree

$$\begin{aligned} IG(\text{Play_basketball}, \text{outlook}) \\ &= H(\text{Play_basketball}) - H(\text{Play_basketball}, \text{outlook}) \\ &= 0.2831 - 0.2087 = 0.0744 \end{aligned}$$

ID3 (Iterative Dichotomiser 3)

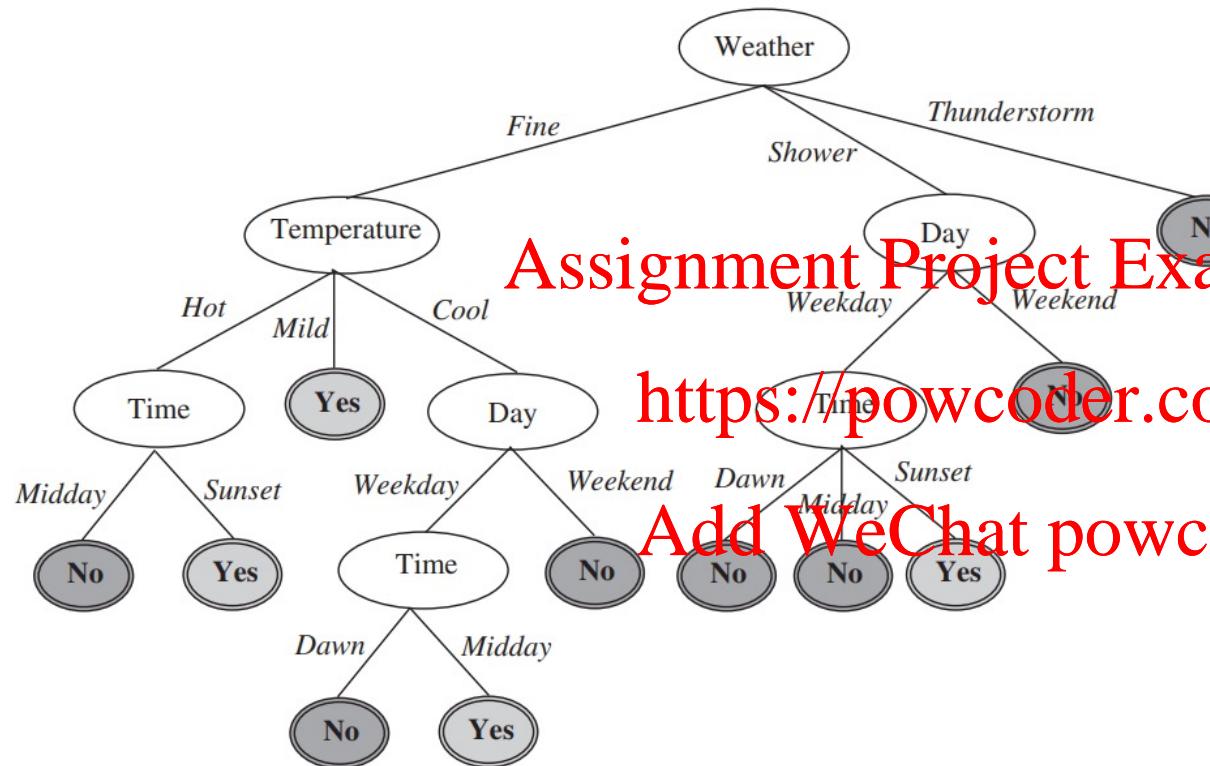
- ❑ It constructs DT, by finding for each node attribute that returns the highest information gain to split the data

Steps

1. Compute the entropy for dataset S $\rightarrow H(S)$
2. For every attribute/feature A:
 - 2.1. Calculate entropy for each categorical value of A $\rightarrow H(S_i)$
 - 2.2. Take weighted average entropy for the current attribute $\rightarrow H(S, A) = \sum_{i \in Values(A)} p_i H(S_i)$
 - 2.3. Calculate IG for the current attribute $\rightarrow IG(S, A) = H(S) - H(S, A)$
3. Pick the attribute with highest IG to be a node, and split dataset by its branch to child nodes/subsets
4. Repeat same process at every child node until the tree is complete

Stopping condition: when data in each partition have same target class

Decision Trees: To Jog or Not To Jog



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

A decision tree is constructed based only on the given training dataset. It is not based on a universal belief.

Figure 17.10 A decision tree

ID3

Example: Consider a piece of data collected over the course of 15 days where the features are Weather, Temperature, Time, Day and the outcome variable is whether Jogging was done on the day. Now, our job is to build a predictive model which takes in above 4 parameters and predicts whether Jogging will be done on the day. We'll build a decision tree to do that using **ID3 algorithm.**

Rec#	Weather	Temperature	Time	Day	Jog (<i>Target Class</i>)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

ID3 (Iterative Dichotomiser 3)

- ❑ It constructs DT, by finding for each node attribute that returns the highest information gain to split the data

Steps

1. Compute the entropy for dataset S

$$H(S)$$

2. For every attribute/feature A :

2.1. Calculate entropy for each categorical value of A

$$H(S_i)$$

2.2. Take weighted average entropy for the current attribute

$$H(S, A) = \sum_{i \in Values(A)} p_i H(S_i)$$

2.3. Calculate IG for the current attribute

$$IG(S, A) = H(S) - H(S, A)$$

- Pick the attribute with highest IG to be a node, and split dataset by its branch to child nodes/subsets
- Repeat same process at every child node until the tree is complete

Stopping condition: when data in each partition have same target class

ID3

Entropy for the given probability of the target classes, p_1, p_2, \dots, p_n where

$\sum_{i=1}^n p_i = 1$, can be calculated as follows:

$$\text{entropy}(p_1, p_2, \dots, p_n) = -\sum_{i=1}^n (p_i \log(1/p_i)) \quad (17.2)$$

$$\begin{aligned} \text{entropy}(\text{Yes}, \text{No}) &= 5/15 \times \log(15/5) + 10/15 \times \log(15/10) \\ &= 0.2764 \end{aligned} \quad (17.3)$$

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

Jog	
Yes	No
5	10

ID3 (Iterative Dichotomiser 3)

- ❑ It constructs DT, by finding for each node attribute that returns the highest information gain to split the data

Steps

1. Compute the entropy for dataset S → $H(S)$
2. For every attribute/feature A:
 - 2.1. Calculate entropy for each categorical value of A → $H(S_i)$
 - 2.2. Take weighted average entropy for the current attribute → $H(S, A) = \sum_{i \in Values(A)} p_i H(S_i)$
 - 2.3. Calculate IG for the current attribute → $IG(S, A) = H(S) - H(S, A)$
3. Pick the attribute with highest IG to be a node, and split dataset by its branch to child nodes/subsets
4. Repeat same process at every child node until the tree is complete

Stopping condition: when data in each partition have same target class

ID3

$$\begin{aligned} \text{entropy}(\text{Weather}=\text{Fine}) &= 4/7 \times \log(7/4) + 3/7 \times \log(7/3) \\ &= 0.2966 \end{aligned} \quad (17.4)$$

$$\begin{aligned} \text{entropy}(\text{Weather}=\text{Shower}) &= 1/4 \times \log(4/1) + 3/4 \times \log(4/3) \\ &= 0.2442 \end{aligned} \quad (17.5)$$

Assignment Project Exam Help

- Step 2: Process attribute *Weather*

<https://powcoder.com>

- Calculate weighted sum entropy of attribute *Weather*:

$$\text{entropy}(\text{Fine}) = 0.2966$$

$$\text{entropy}(\text{Shower}) = 0.2442$$

$$\text{entropy}(\text{Thunderstorm}) = 0 + 4/4 \times \log(4/4) = 0$$

$$\text{weighted sum entropy}(\text{Weather}) = 0.2035$$

- Calculate information gain for attribute *Weather*:

$$\text{gain}(\text{Weather}) = 0.0729$$

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

		Jog		
		Yes	No	
Weather	Fine	4	3	7
	Shower	1	3	4
	Thunderstorm	0	4	4
				15

ID3

Weighted sum entropy (*Weather*) = Weighted entropy (*Fine*)
 + Weighted entropy (*Shower*)
 + Weighted entropy (*Thunderstorm*)
 $= 7/15 \times 0.2966 + 4/15 \times 0.2442 + 4/15 \times 0$
 $= 0.2035$

Assignment Project Exam Help
(17.8)

- Step 2: Process attribute *Weather*

<https://powcoder.com>

- Calculate weighted sum entropy of attribute *Weather*:

$$\text{entropy}(\text{Fine}) = 0.2966$$

$$\text{entropy}(\text{Shower}) = 0.2442$$

$$\text{entropy}(\text{Thunderstorm}) = 0 + 4/4 \times \log(4/4) = 0$$

$$\text{weighted sum entropy}(\text{Weather}) = 0.2035$$

- Calculate information gain for attribute *Weather*:

$$\text{gain}(\text{Weather}) = 0.0729$$

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

		Jog		
		Yes	No	
Weather	Fine	4	3	7
	Shower	1	3	4
	Thunderstorm	0	4	4
				15

ID3

$$\begin{aligned} \text{gain}(Weather) &= \text{entropy}(\text{training dataset } D) - \text{entropy}(\text{attribute } Weather) \\ &= 0.2764 - 0.2035 \\ &= 0.0729 \end{aligned} \tag{17.7}$$

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

- Step 2: Process attribute *Weather*

<https://powcoder.com>

- Calculate weighted sum entropy of attribute *Weather*:

$$\text{entropy}(\text{Fine}) = 0.2966$$

Add WeChat ^(equation 17.4) ~~powcoder~~ ^(equation 17.4)

$$\text{entropy}(\text{Shower}) = 0.2442$$

$$\text{entropy}(\text{Thunderstorm}) = 0 + 4/4 \times \log(4/4) = 0$$

$$\text{weighted sum entropy}(\text{Weather}) = 0.2035 \tag{equation 17.6}$$

- Calculate information gain for attribute *Weather*:

$$\text{gain}(\text{Weather}) = 0.0729$$

(equation 17.7)

ID3

- Step 3: Process attribute *Temperature*

- Calculate weighted sum entropy of attribute *Temperature*:

$$\text{entropy}(\text{Hot}) = 2/5 \times \log(5/2) + 3/5 \times \log(3/1) = 0.2923$$

$$\text{entropy}(\text{Mild}) = \text{entropy}(\text{Hot})$$

$$\text{entropy}(\text{Cool}) = 1/5 \times \log(5/1) + 4/5 \times \log(5/4) = 0.2173$$

$$\text{weighted sum entropy}(\text{Temperature}) = 5/15 \times 0.2923 + 5/15 \times 0.2173 \\ = 0.2674$$

- Calculate information gain for attribute *Temperature*:

$$\text{gain}(\text{Temperature}) = 0.2764 - 0.2674 = 0.009$$

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

		Jog		
		Yes	No	
Temperature	Hot	2	3	5
	Mild	3	2	5
	Cool	1	4	5
				15

ID3

- Step 4: Process attribute *Time*

Assignment Project Exam Help

- Calculate weighted sum entropy of attribute *Time*:

$$\text{entropy}(\text{Dawn}) = 0 + 5/5 \times \log(1/5) = 0$$

$$\text{entropy}(\text{Midday}) = 2/6 \times \log(6/2) + 4/6 \times \log(6/4) = 0.2764$$

$$\text{entropy}(\text{Sunset}) = 3/4 \times \log(4/2) + 1/4 \times \log(4/1) = 0.2443$$

$$\text{weighted sum entropy } (\text{Time}) = 0 + 6/15 \times 0.2764 + 4/15 \times 0.2443 =$$

$$0.1757$$

- Calculate information gain for attribute *Time*:

$$\text{gain } (\text{Temperature}) = 0.2764 - 0.1757 = 0.1007$$

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

		Jog		
		Yes	No	
Time	Dawn	0	5	5
	Midday	2	4	6
	Sunset	3	1	4
				15

ID3

Assignment Project Exam Help

Step 5: Process attribute *Day*

- Calculate weighted sum entropy of attribute *Day*:

$$\text{entropy}(\text{Weekday}) = 4/10 \times \log(10/4) + 6/10 \times \log(10/6) \\ = 0.2923$$

$$\text{entropy}(\text{Weekend}) = 1/5 \times \log(1/1) + 4/5 \times \log(5/4) \\ = 0.2173$$

$$\text{weighted sum entropy (Day)} = 10/15 \times 0.2923 + 5/15 \\ \times 0.2173 = 0.2674$$

- Calculate information gain for attribute *Day*:

$$\text{gain}(\text{Temperature}) = 0.2764 - 0.2674 = 0.009$$

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

		Jog	
		Yes	No
Day	Weekend	4	6
	Weekday	1	4
			15

ID3 (Iterative Dichotomiser 3)

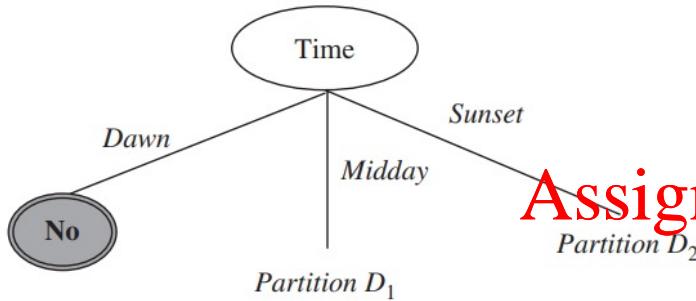
- ❑ It constructs DT, by finding for each node attribute that returns the highest information gain to split the data

Steps

1. Compute the entropy for dataset S → $H(S)$
2. For every attribute/feature A:
 - 2.1. Calculate entropy for each categorical value of A → $H(S_i)$
 - 2.2. Take weighted average entropy for the current attribute → $H(S, A) = \sum_{i \in Values(A)} p_i H(S_i)$
 - 2.3. Calculate IG for the current attribute → $IG(S, A) = H(S) - H(S, A)$
3. Pick the attribute with highest IG to be a node, and split dataset by its branch to child nodes/subsets
4. Repeat same process at every child node until the tree is complete

Stopping condition: when data in each partition have same target class

ID3



Assignment Project Exam Help

Figure 17.13 Attribute *Time* as the root node

<https://powcoder.com>

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

Comparing equations 17.7, 17.8, 17.9, and 17.10 ,and 17.10 for the gain of each other attributes (Weather, Temperature, Time, and Day), the biggest gain is *Time*, with gain value = 0.1007 (see equation 17.9), and as a result, attribute *Time* is chosen as the first splitting attribute. A partial decision tree with the root node *Time* is shown in Figure 17.13.

ID3

Jog	
Yes	No
2	4

		Jog		
		Yes	No	
Day	Weekend	0	0	0
	Weekday	2	4	6
				6
		Yes	No	
Weather	Fine	2	1	3
	Shower	0	1	1
	Thunderstorm	0	2	2
				6
		Yes	No	
Temperature	Hot	0	2	2
	Mild	1	1	2
	Cool	1	1	2
				6

Assignment Project Exam Help

- The next stage is to process partition D_1 consisting of records with Time=Midday. Training dataset partition D_1 consists of 6 records with record#: 3, 6, 8, 9, 10, and 15. The next task is to determine the splitting attribute for partition D_1 , whether it is *Weather*, *Temperature*, or *Day*.

Step 1: Calculate entropy for the training dataset partition D_1 .

$$\text{entropy}(D_1) = 2/6 \log(6/2) + 4/6 \log(6/4) = 0.2764 \quad (17.11)$$

Step 2: Process attribute *Weather*

- Calculate weighted sum entropy of attribute *Weather*

$$\text{entropy}(\text{Fine}) = 2/3 \times \log(6/2) + 1/3 \times \log(3/1) = 0.2764$$

$$\text{entropy}(\text{Shower}) = \text{entropy}(\text{Thunderstorm}) = 0$$

$$\text{weighted sum entropy}(\text{Weather}) = 3/5 \times 0.2764 = 0.1382$$

- Calculate information gain for attribute *Weather*:

$$\text{gain}(\text{Weather}) = 0.2764 - 0.1382 = 0.1382 \quad (17.12)$$

Step 3: Process attribute *Temperature*

- Calculate weighted sum entropy of attribute *Temperature*

$$\text{entropy}(\text{Hot}) = 0$$

$$\text{entropy}(\text{Mild}) = \text{entropy}(\text{Cool}) = 1/2 \times \log(2/1) + 1/2 \times \log(2/1) = 0.3010$$

$$\text{weighted sum entropy}(\text{Temperature}) = 2/6 \times 0.3010 + 2/6 \times 0.3010 = 0.2006$$

- Calculate information gain for attribute *Temperature*:

$$\text{gain}(\text{Temperature}) = 0.2764 - 0.2006 = 0.0758 \quad (17.13)$$

Step 4: Process attribute *Day*

- Calculate weighted sum entropy of attribute *Day*:

$$\text{entropy}(\text{Weekday}) = 2/6 \times \log(6/2) + 4/6 \times \log(6/4) = 0.2764$$

$$\text{entropy}(\text{Weekend}) = 0$$

$$\text{weighted sum entropy}(\text{Day}) = 0.2764$$

- Calculate information gain for attribute *Day*:

$$\text{gain}(\text{Temperature}) = 0.2764 - 0.2764 = 0 \quad (17.14)$$

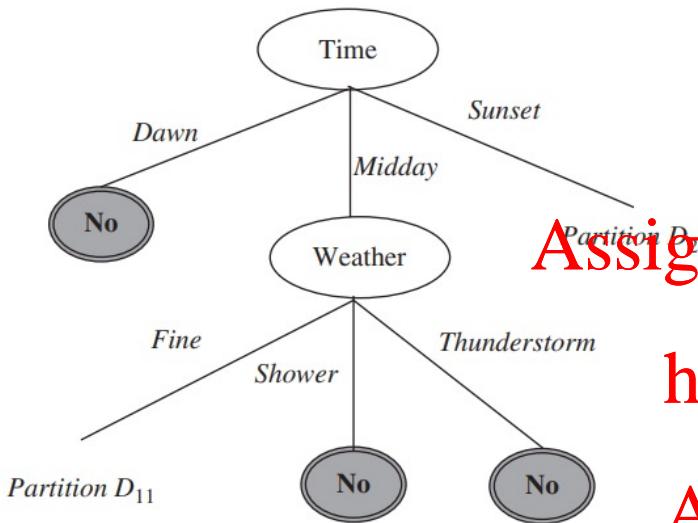
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The best splitting node for partition D_1 is attribute **Weather** with information gain value of 0.1382 (see equation 17.12).

ID3



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- The next stage is to process partition D_1 consisting of records with Time=Midday. Training dataset partition D_1 consists of 6 records with record#: 3, 6, 8, 9, 10, and 15. The next task is to determine the splitting attribute for partition D_1 , whether it is Weather, Temperature, or Day.

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

ID3

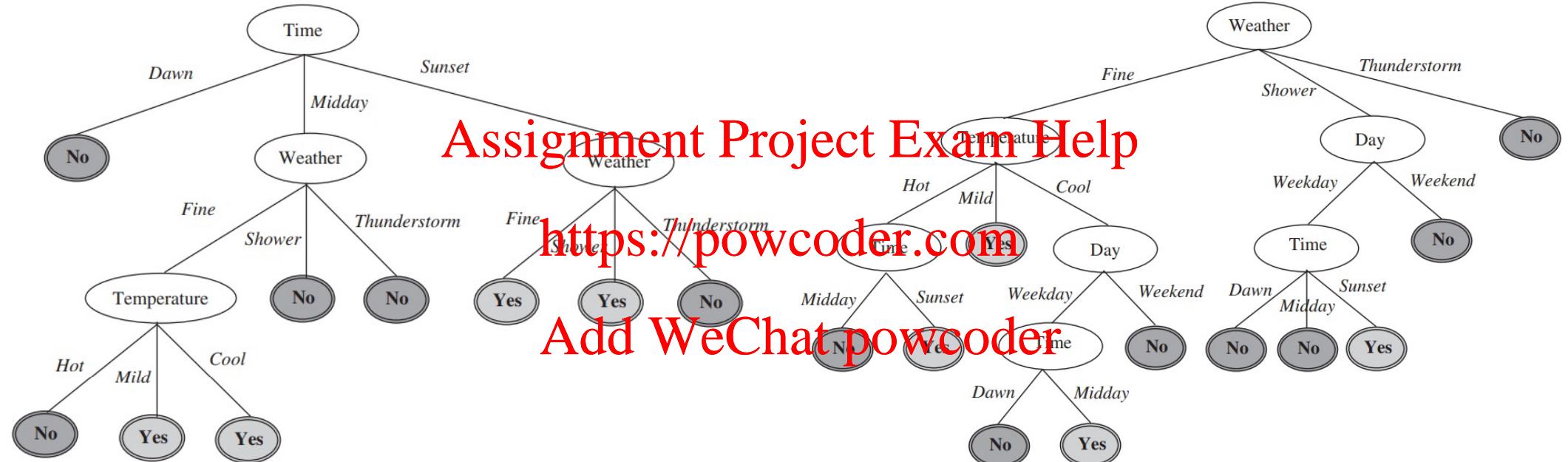


Figure 17.15 Final decision tree

Figure 17.10 A decision tree

Maximum Depth of DT

maxDepth: the largest possible length between the root to a leaf (or maximum level of the tree).



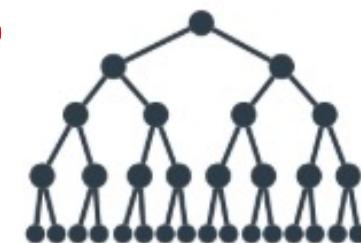
Depth = 1



Depth = 2



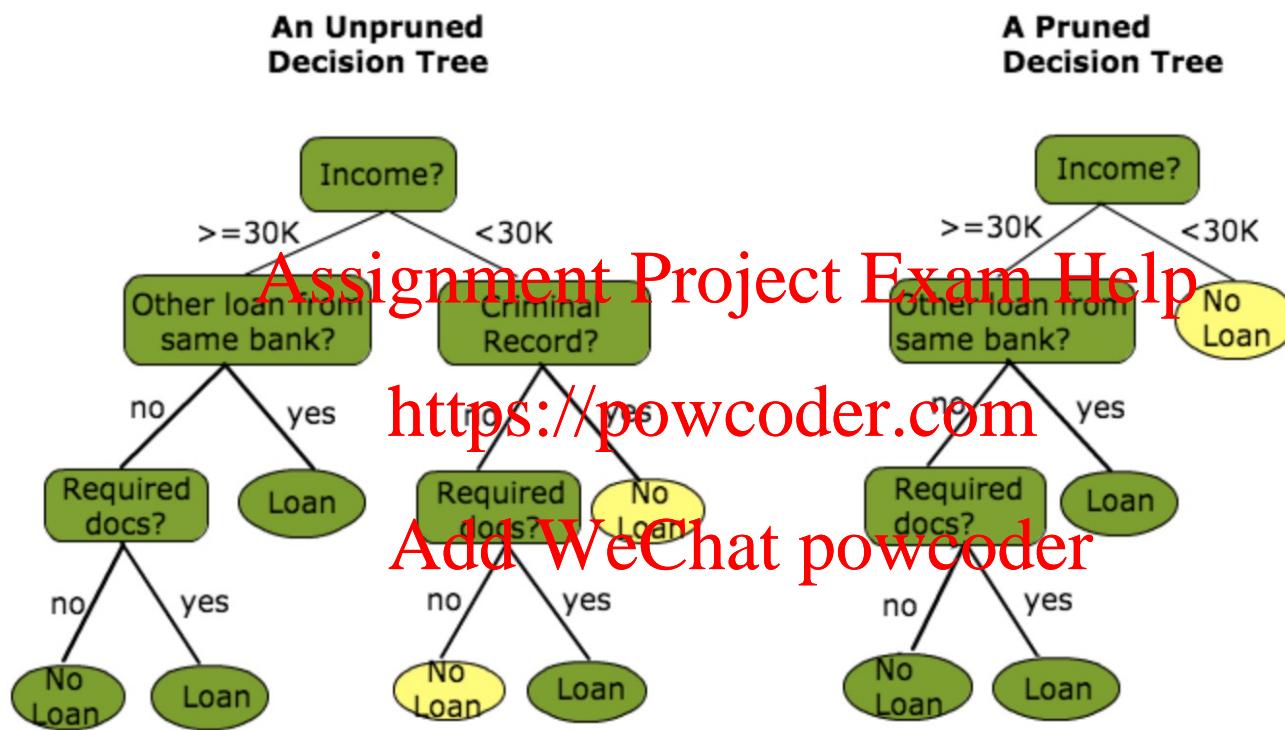
Depth = 3



Depth = 4

Maximum depth of a decision tree

Pruning



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

<https://kaumadiechamalka100.medium.com/decision-tree-in-machine-learning-c610ef087260>

Decision Tree Algorithm

Advantages:

- Easy to understand.
- Easy to generate rules.
- There are almost null hyper-parameters to be tuned.
- Complex Decision Tree models can be significantly simplified by its visualizations.

Assignment Project Exam Help

<https://powcoder.com>

Disadvantages:

- Might suffer from overfitting.
- Does not easily work with non-numerical data.
- Low prediction accuracy for a dataset in comparison with other machine learning classification algorithms.
- When there are many class labels, calculations can be complex.

Add WeChat powcoder

Ensemble methods

A single decision tree have the tendency to overfit

But, it is super fast.

Assignment Project Exam Help

How about multiple trees at once?

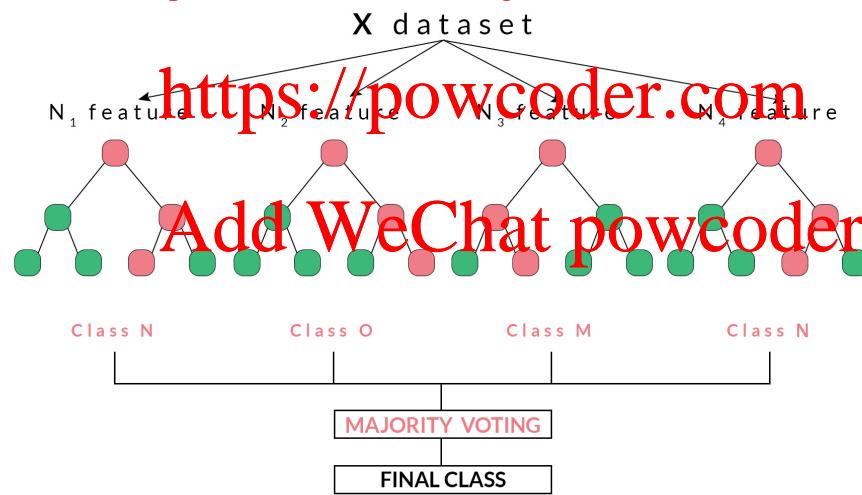
<https://powcoder.com>

Add WeChat powcoder
Make sure they do not all just learn the same!

Random Forest Algorithm

Random forest (or **random forests**) is an ensemble classifier that consists of many decision trees and outputs the class that is the mode of the class's output by individual trees.

Assignment Project Exam Help



Optimisations

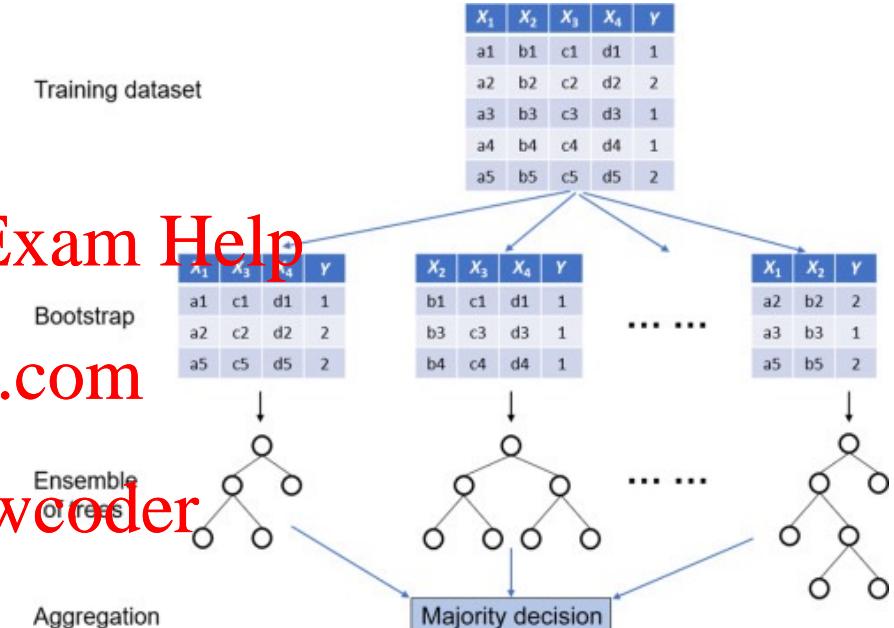
1. Bagging: Bootstrap **aggregating** is a method that result in low variance – used to reduce variance of DTs

Assignment Project Exam Help

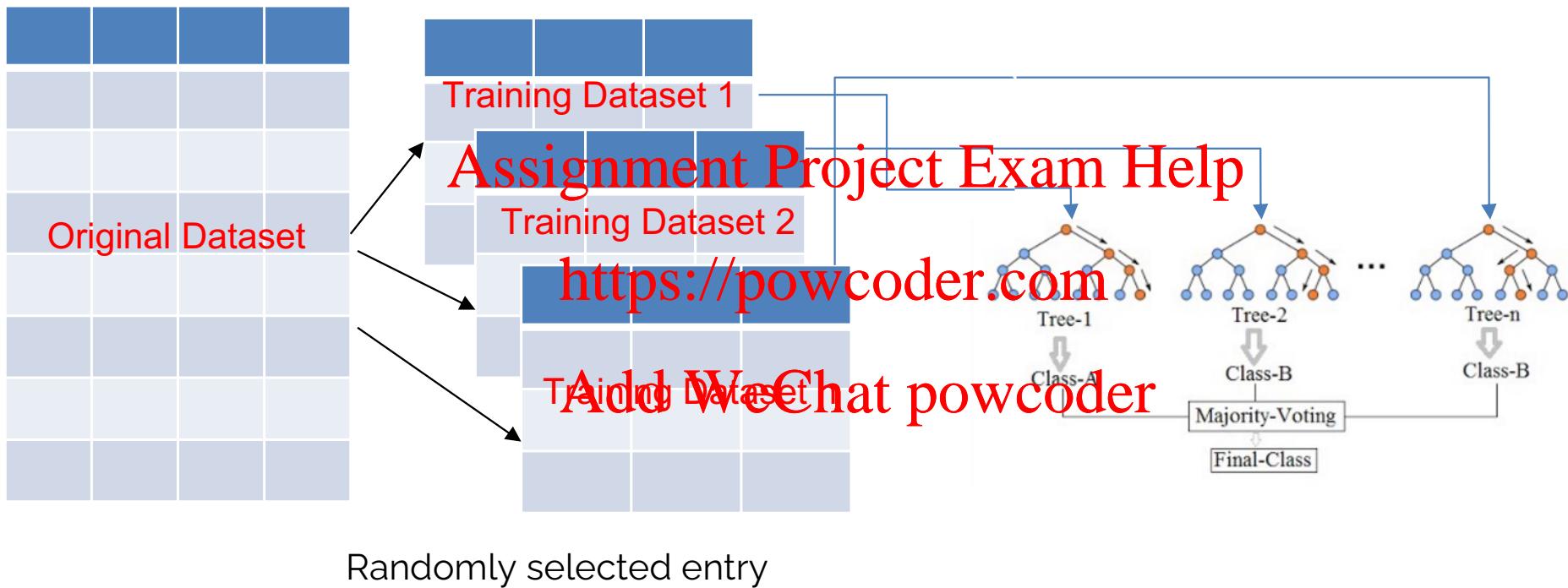
Rather than training each tree on all the inputs in the training set (producing multiple identical trees), each tree is trained on different set of sample data

<https://powcoder.com>

Add WeChat powcoder



Example



Optimisations

2. Gradient boosting: selecting best classifiers to improve prediction accuracy with each new tree.

Assignment Project Exam Help

- It works by combining several weak learners (typically high bias, low variance models) to produce an overall strong model.
<https://powcoder.com>
- It builds one tree at a time, works in a forward stage-wise manner, - adding a classifier at a time, so that the next classifier is trained to improve the already trained ensemble.
[Add WeChat powcoder](#)

Advantages and Disadvantages of Random Forest

Advantages

- It is robust to correlated predictors.
- It is used to solve both regression and classification problems.
- It can be also used to solve unsupervised ML problems.
- It can handle thousands of input variables without variable selection.
- It can be used as a feature selection tool using its variable importance plot.
- It takes care of missing data internally in an effective manner.

Advantages and Disadvantages of Random Forest

Disadvantages

- The Random Forest model is difficult to interpret.
- It tends to return erratic predictions for observations out of range of training data. For example, the training data contains two variable x and y. The range of x variable is 30 to 70. If the test data has x = 200, random forest would give an unreliable prediction.
- It can take longer than expected time to computer a large number of trees.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

What have we learnt today?

Classification techniques

Decision Trees and Random Forest and KNN

When and how to use each technique

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help Parallel Classification

<https://powcoder.com>

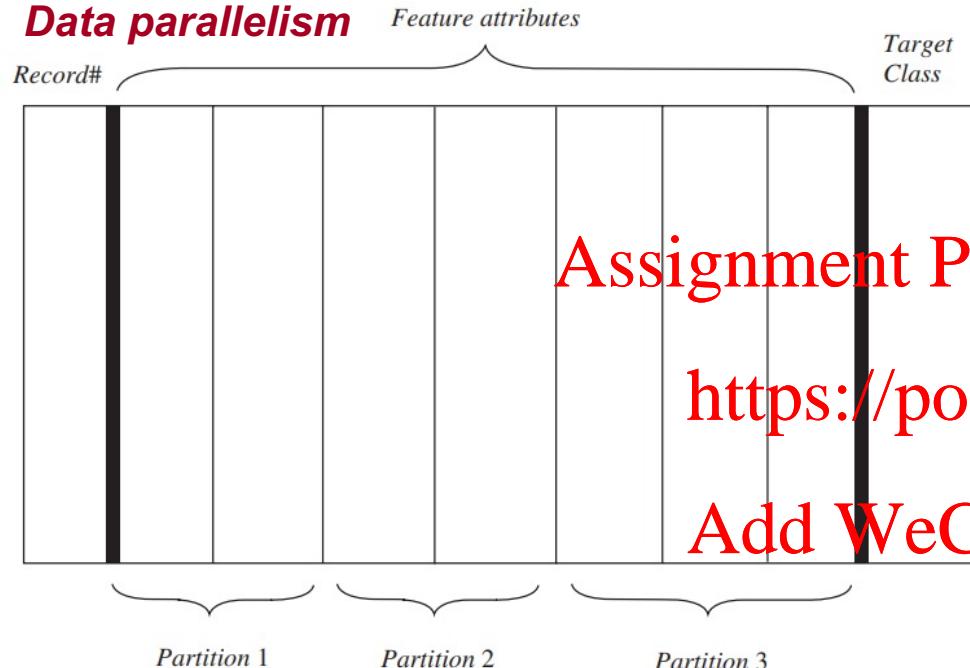
Add WeChat powcoder

Prajwol Sangat



Parallel Classification: Decision Tree

Data parallelism



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

Figure 17.16 Vertical data partitioning of training data set

Parallel Classification: Decision Tree

Data parallelism: Vertical Partitioning of Training dataset

Rec#	Weather	Temperature	Jog (Target Class)
1	Fine	Mild	Yes
2	Fine	Hot	Yes
3	Shower	Mild	No
4	Thunderstorm	Cool	No
5	Shower	Hot	Yes
6	Fine	Hot	No
7	Fine	Cool	No
8	Thunderstorm	Cool	No
9	Fine	Cool	Yes
10	Fine	Mild	No
11	Shower	Hot	No
12	Shower	Mild	No
13	Fine	Cool	No
14	Thunderstorm	Mild	No
15	Thunderstorm	Hot	No

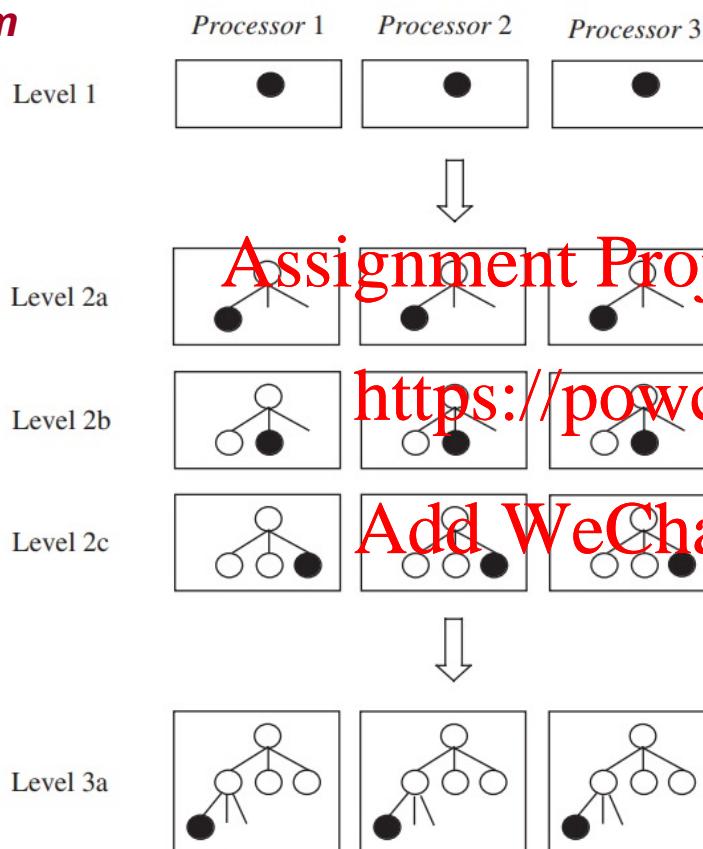
Partition 1

Rec#	Time	Day	Jog (Target Class)
1	Sunset	Weekend	Yes
2	Sunset	Weekday	Yes
3	Midday	Weekday	No
4	Dawn	Weekend	No
5	Sunset	Weekday	Yes
6	Midday	Weekday	No
7	Dawn	Weekend	No
8	Midday	Weekday	No
9	Midday	Weekday	Yes
10	Midday	Weekday	Yes
11	Dawn	Weekend	No
12	Dawn	Weekday	No
13	Dawn	Weekday	No
14	Sunset	Weekend	No
15	Midday	Weekday	No

Partition 2

Parallel Classification: Decision Tree

Data parallelism



Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

Figure 17.17 Data parallelism of parallel decision tree construction

Parallel Classification: Decision Tree

Level 1 (Root Node):

Data parallelism

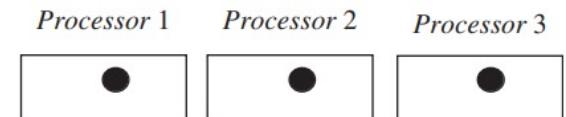
Processor 1

Rec#	Weather	Temperature	Target Class
1			
2			
...			
15			

Processor 2

Rec#	Time	Day	Target Class
1			
2			
...			
15			

Level 1

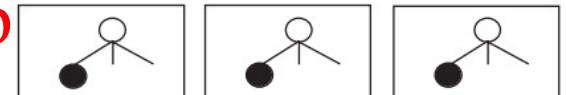


Locally calculate the information gain values for: *Weather* and *Temperature*

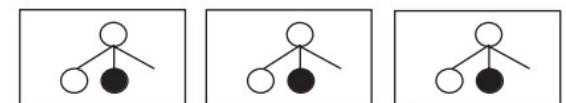


Locally calculate the information gain values for: *Time* and *Day*

Level 2a



Level 2b



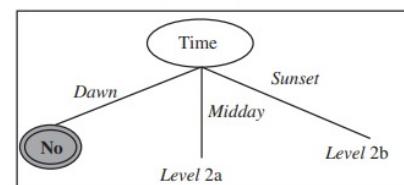
Global information sharing stage:

- Share target class counts to calculate dataset entropy value
- Exchange dataset entropy value to determine splitting attribute (e.g. Time attribute is decided to be the splitting attribute)
- Distribute selected records# to all processor for the next phase (e.g. records 3, 6, 8, 9, 10, 15 for Time *Midday*, and records 1, 2, 5, 14 for Time *Sunset*)

Add WeChat powcoder

Decision tree for Level 1:

Processor 1



Processor 2

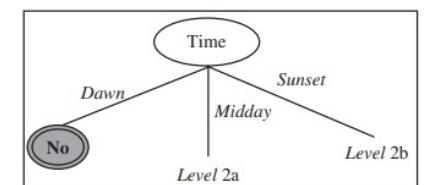


Figure 17.18 Data parallelism in decision tree

Parallel Classification: Decision Tree

Data parallelism

Level 2a:

Processor 1	Rec#	Weather	Temperature	TargetClass
	3			
	6			
	8			
	9			
	10			
	15			

Processor 2	Rec#	Time	Day	TargetClass
	3			
	6			
	9			
	10			
	15			

Locally calculate the information gain values for: Weather and Temperature

Locally calculate the information gain values for Day

Assignment Project Exam Help

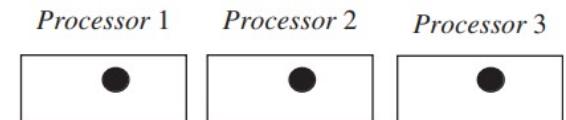
<https://powcoder.com>

Add WeChat powcoder

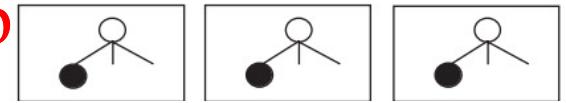
Global information sharing stage:

- Share target class counts of each partition to calculate dataset entropy value
- Exchange dataset entropy value to determine splitting attribute (e.g. Weather attribute is decided to be the splitting attribute)
- Distribute selected records# to all processor for the next phase

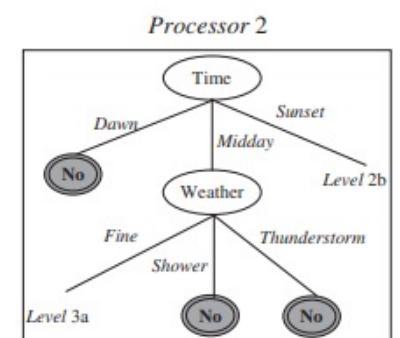
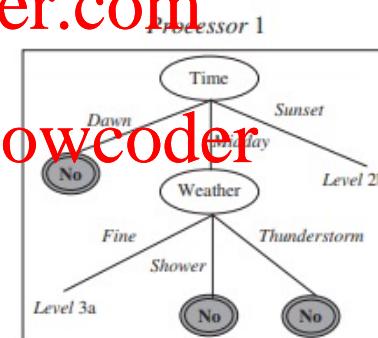
Level 1



Level 2a



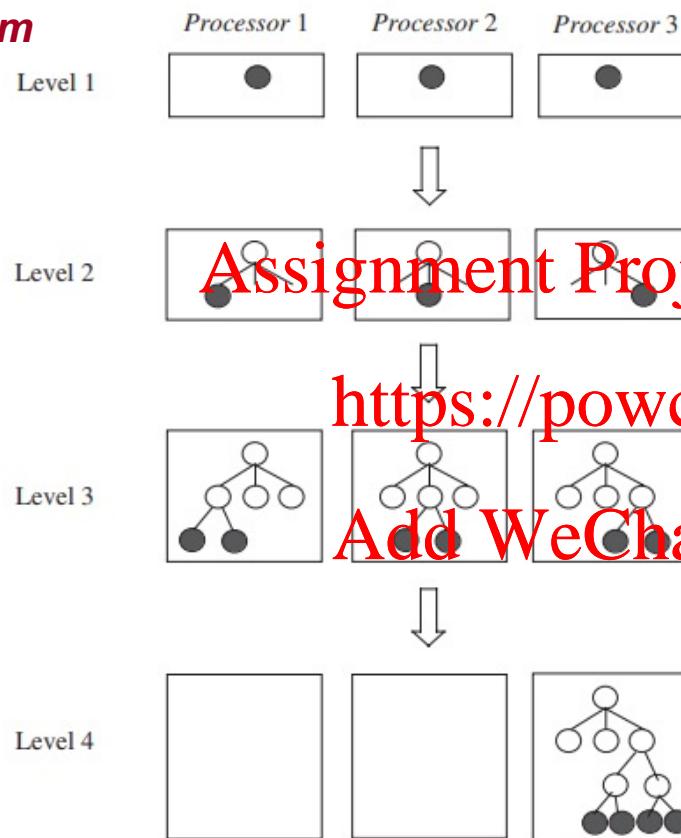
Result decision tree for Level 2:



Level 2b: to continue...

Parallel Classification: Decision Tree

Result parallelism



Rec#	Weather	Temperature	Time	Day	Jog (Target Class)
1	Fine	Mild	Sunset	Weekend	Yes
2	Fine	Hot	Sunset	Weekday	Yes
3	Shower	Mild	Midday	Weekday	No
4	Thunderstorm	Cool	Dawn	Weekend	No
5	Shower	Hot	Sunset	Weekday	Yes
6	Fine	Hot	Midday	Weekday	No
7	Fine	Cool	Dawn	Weekend	No
8	Thunderstorm	Cool	Midday	Weekday	No
9	Fine	Cool	Midday	Weekday	Yes
10	Fine	Mild	Midday	Weekday	Yes
11	Shower	Hot	Dawn	Weekend	No
12	Shower	Mild	Dawn	Weekday	No
13	Fine	Cool	Dawn	Weekday	No
14	Thunderstorm	Mild	Sunset	Weekend	No
15	Thunderstorm	Hot	Midday	Weekday	No

Figure 17.11. Training dataset

Figure 17.20 Result parallelism of parallel decision tree construction

Parallel Classification: Decision Tree

Result parallelism

Horizontal Data Partitioning:

Processor 1						Processor 2					
Rec#	Weather	Temp	Time	Day	Target Class	Rec#	Weather	Temp	Time	Day	Target Class
1						9					
2						10					
...						...					
8						15					

Assignment Project Exam Help

<https://powcoder.com>

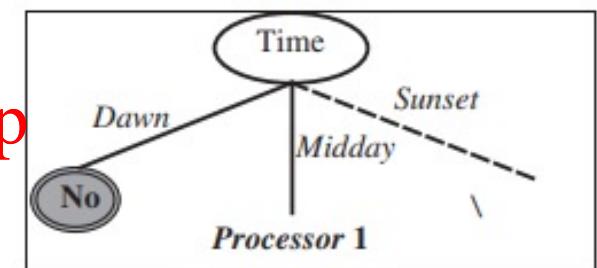
Level 1 (Root Node):

- Count target class on each partition
- Perform intra-node parallelism the same as for data parallelism to share target class counts to calculate dataset entropy value, exchange dataset entropy value to determine splitting attribute, and distribute selected records# to all other processors for the next phase)

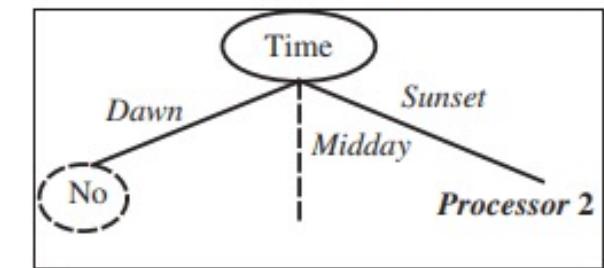
Add WeChat powcoder

Decision tree for Level 1:

Processor 1



Processor 2



Parallel Classification: Decision Tree

Result parallelism

Jog	
Yes	No
5	10

		Jog		
		Yes	No	
Time	Dawn	0	5	5
	Midday	2	4	6
	Sunset	3	7	4

Assignment Project Exam Help

		Jog		
		Yes	No	
Day	Weekend	4	6	10
	Weekday	1	4	5

<https://powcoder.com>

Add WeChat powcoder

		Jog		
		Yes	No	
Weather	Fine	4	3	7
	Shower	1	3	4
	Thunderstorm	0	4	4

9

Parallel Classification: Decision Tree

Result parallelism

Level 2:

Processor 1						Processor 2					
Rec#	Weather	Temp	Time	Day	Target Class	Rec#	Weather	Temp	Time	Day	Target Class
3						2					
6						5					
8						14					
9											
10											
15											

Assignment Project Exam Help
<https://powcoder.com>



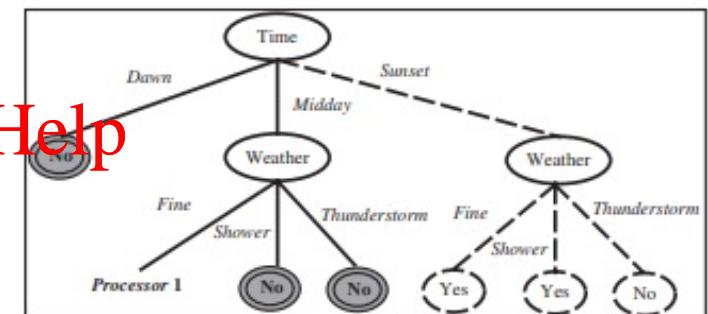
Add WeChat powcoder

Global information sharing stage:

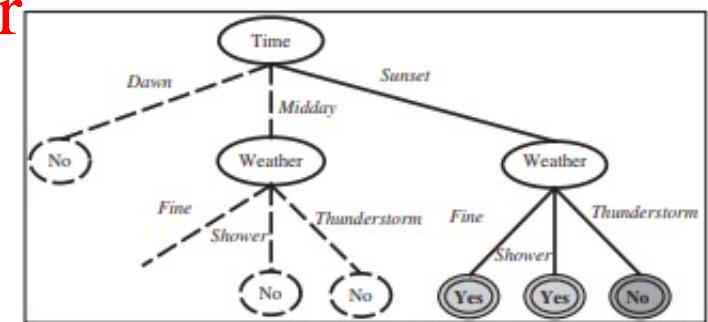
- Count target class on each partition
- Perform intra-node parallelism the same as for data parallelism to share target class counts to calculate dataset entropy value, exchange dataset entropy value to determine splitting attribute, and distribute selected records# to all other processors for the next phase)

Result decision tree for Level 2:

Processor 1



Processor 2



Parallel Classification: Decision Tree

Level 3:

Processor 1	Rec#	Weather	Temp	Time	Day	Target Class
	6					
	9					
	10					

Processor 2	Rec#	Weather	Temp	Time	Day	Target Class

Assignment Project Exam Help

<https://powcoder.com>

Global information sharing stage:... as like in Level 1

Result decision tree for Level 3:

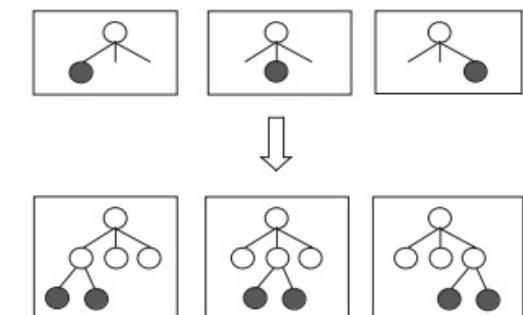
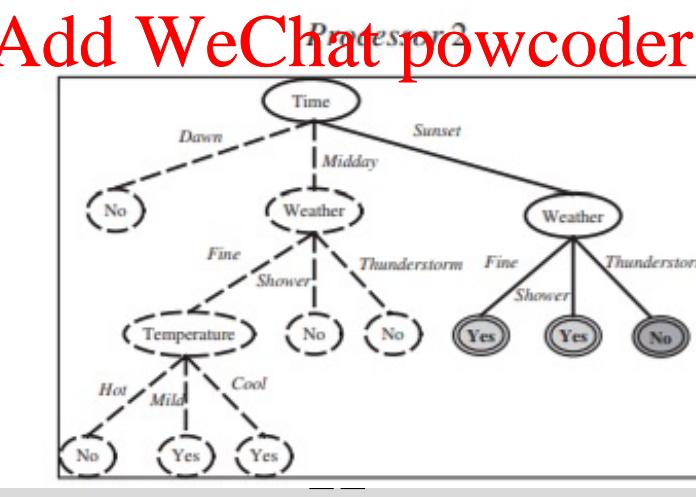
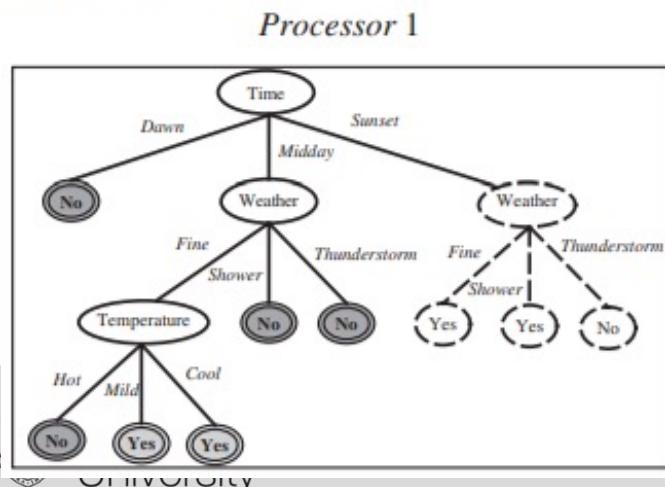


Figure 17.20 Result parallelism of parallel decision tree construction

See you next week..

**Assignment Project Exam Help
Questions???**

<https://powcoder.com>

Add WeChat powcoder