

# Lecture Scribe For Machine Learning(CS60050),

Spring 2020-2021

## Decision Tree

Date: 14th January, 2021

The screenshot shows a Zoom meeting window with a hand-drawn lecture on Decision Tree Learning. The lecture is written on a green chalkboard background and includes the following content:

- Boolean AND formula:**  $h_1 = (\text{Sky} = \text{Sunny} \wedge \text{Water} = \text{Warm}) \vee (\text{Sky} = \text{Cloudy}) \vee (\text{Sky} = \text{Raining} \wedge \text{Humidity} = \text{High})$
- Learning Process Diagram:** A circular diagram labeled "Learning" with "TE" and "S" inputs, and "What?" and "Decision Tree Learning" outputs.
- Decision Tree:** A tree structure for weather classification. The root node is "Sky" (Sunny, Cloudy, Raining). "Sunny" leads to "Water" (Warm, Cool). "Warm" leads to "Yes (+)", "Cool" leads to "No (-)". "Cloudy" leads to "Yes". "Raining" leads to "Humidity" (High, Normal). "High" leads to "Yes (+)", "Normal" leads to "No (-)".
- Decision Tree Definition:**  $\mathcal{H} = \{ \text{Decision Tree} \}$  with variables  $\{v_1, v_2, v_3\}$  and formula  $(a \wedge b) \vee (c \wedge d) \vee (a \wedge c \wedge e)$  where  $a \neq v_1$ .

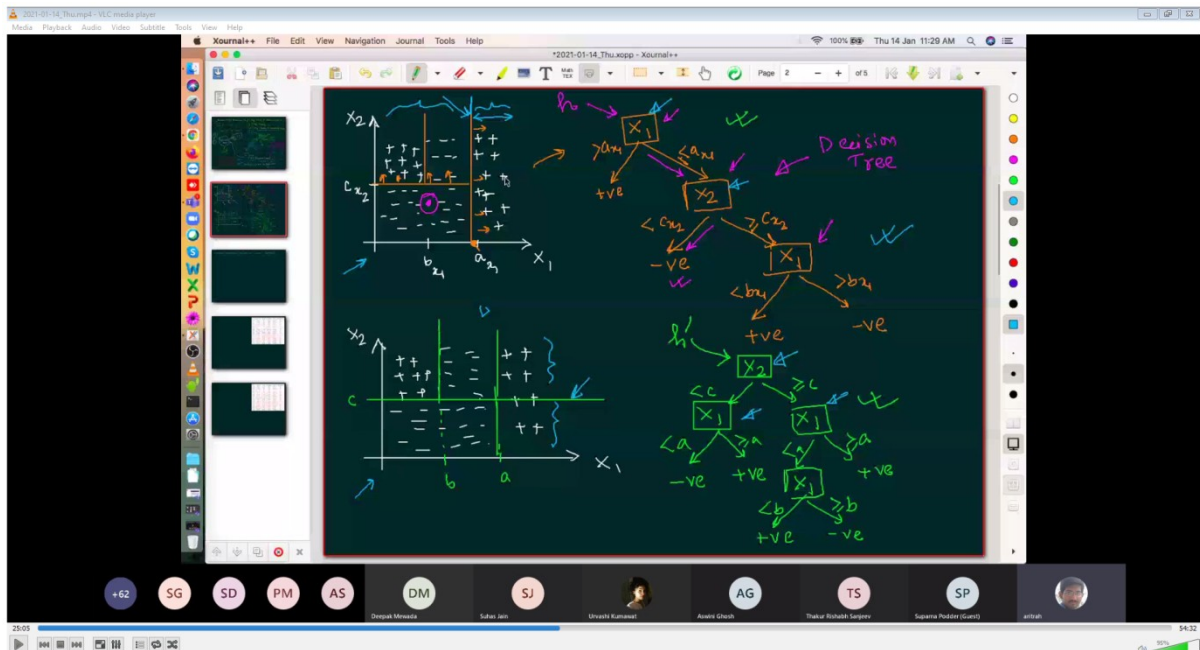
The Zoom interface at the bottom shows a meeting ID of 1401, a time of 54:32, and several participants: SG, SD, PM, RT, DM, SJ, AG, AS, SP.

## Decision Tree Learning

Decision tree learning is a method for approximating discrete-valued target functions. The learned function is represented by a decision tree. Decision tree learning is one of the most widely used and practical methods for inductive inference. Decision tree learning method searches a completely expressive hypothesis.

# Decision Tree

Decision trees represent a disjunction of conjunctions of constraints on the attribute values of instances. Each path from the tree root to a leaf corresponds to a conjunction of attribute tests. The tree itself is a disjunction of these conjunctions. Decision trees classify instances by sorting them down the tree from the root to some leaf node, which provides the classification of the instance. Each node in the tree specifies a test of some attribute of the instance. Each branch descending from a node corresponds to one of the possible values for the attribute. Each leaf node assigns a classification.



## Which Attribute is most suitable

We would like to select the attribute that is most useful for classifying Examples. **Information gain** measures how well a given attribute separates the training examples according to their target classification. In order to define information gain precisely, we use a measure commonly used in information theory, called **entropy**.

## Entropy

The whiteboard content includes:

- Question: "Which attribute should we choose first/next?"
- Section: "Entropy"
  - Formula:  $E = -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$
  - Diagram: A decision tree with root node  $X_1$  and child nodes  $[x_1, x_2]$ ,  $[x_3, x_4]$ , and  $[x_5, x_6]$ . Weights  $w_1, w_2, w_3$  are shown.
  - Formula:  $E = -\sum_{c=1}^m p_c \log_2 p_c$
- Section: "Information Gain"
  - Formula:  $IG(S, A) = E(S) - \sum_{v=1}^{|A|} \frac{|S_v|}{|S|} Entropy(S_v)$

Entropy characterizes the impurity of an arbitrary collection of examples.

Given a collection  $S$ , containing positive and negative examples of some target concept, the entropy of  $S$  relative to this Boolean classification is:

$$\text{Entropy}(S) = -p^+ \log_2 p^+ - p^- \log_2 p^-$$

- $S$  is a sample of training examples
- $p^+$  is the proportion of positive examples
- $p^-$  is the proportion of negative examples

## Entropy – Non-Boolean Target Classification

If the target attribute can take on  $c$  different values, then the entropy of  $S$  relative to this  $c$ -wise classification is defined as

$$\text{Entropy}(S) = \text{summation}(i=1 \text{ to } i=c) -p_i \log_2 p_i$$

- $p_i$  is the proportion of  $S$  belonging to class  $i$ .
- The logarithm is still base 2 because entropy is a measure of the expected encoding length measured in bits.
- If the target attribute can take on  $c$  possible values, the entropy can be as large as  $\log_2 c$ .

## Information Gain

- **entropy** is a measure of the impurity in a collection of training examples
- **information gain** is a measure of the effectiveness of an attribute in classifying the training data.

- **information gain** measures the expected reduction in entropy by partitioning the examples according to an attribute.

$$IG(S,A) = Entropy(S) - \sum_{v=1 \text{ to } v=|A|} ( |S_v| / |S| ) Entropy(S_v)$$

- S - a collection of examples
- A - an attribute
- S<sub>v</sub> - the subset of S for which attribute A has value v

## Gini Impurity

Gini Impurity is a measurement used to build Decision Trees to determine how the features of a data set should split nodes to form the tree. More precisely, the Gini Impurity of a data set is a number between 0-0.5, which indicates the likelihood of new, random data being miss classified if it were given a random class label according to the class distribution in the data set.

## Entropy vs Gini Impurity

The maximum value for entropy is 1 whereas the maximum value for Gini impurity is 0.5.

As the Gini Impurity does not contain any logarithmic function to calculate it takes less computational time as compared to entropy.

## Conclusion

we have covered a lot of details about decision tree, how it works and maths behind it, attribute selection measures such as Entropy, Information Gain, Gini Impurity with their formulas, and how machine learning algorithm solves it.

Scribe prepared by Sudhanshu Sharma(18CS30041)