

Decision Tree Algorithm – Regression Problem:

[1]. Read Data, Identify Target Variable and Input Features (Attributes)

[2]. Calculate SD, CV, Mean using below equations

$$\text{Mean} = \frac{\sum x}{n}$$

$$\text{SD} = \text{Square Root of } \frac{\sum (x - \text{mean})^2}{n}$$

$$\text{CV} = \frac{\text{SD}}{\text{mean}}$$

[3]. The dataset is then split on the different attributes. The standard deviation for each branch is calculated.

$$\text{SD}(\text{Attribute}) = \sum_{\text{Branch} \in \text{Attribute}} W(\text{branch}) \text{SD}(\text{Branch})$$

[4]. The resulting standard deviation is subtracted from the standard deviation before the split. The result is the standard deviation reduction.

$$\text{SDR} = \text{SD} - \text{SD}(\text{Attribute})$$

[5]. The attribute with the largest standard deviation reduction is chosen for the decision node.

The dataset is divided based on the values of the selected attribute. This process is run recursively on the non-leaf branches, until all data is processed. (Termination: coefficient of deviation (CV) for a branch becomes smaller than a certain threshold (e.g., 10%) and/or when too few instances (n) remain in the branch)

Decision Tree Algorithm – Classification Problem – Information Gain Index:

[1]. Read Data, Identify Target Variable and Input Features (Attributes)

[2]. Calculate entropy of Target Variable using below equation

$$\text{Entropy}(\text{Target}) = - \sum_{\text{Class} \in \text{Target}} p(\text{Class}) \log_2(p(\text{class}))$$

[3]. Calculate entropy of each attribute using below equation

$$E(Target, Attribute) = \sum_{Branch \in Attribute} p(branch)E(branch)$$

$$where E(branch) = E(a_1, a_2, \dots, a_n) = - \sum_{i \in n} \frac{a_i}{\sum_{i=1}^n a_i} \log_2 \left(\frac{a_i}{\sum_{i=1}^n a_i} \right)$$

[6]. Calculate Information Gain using below equation

$$IG(Target, Attribute) = E(Target) - E(Target, Attribute)$$

[7]. Choose attribute with highest information gain as decision node

[8]. Repeat step 2 to 7 until complete tree formed with leaf node

Decision Tree Algorithm – Classification Problem – Gini Index:

[1]. Read Data, Identify Target Variable and Input Features (Attributes)

[2]. Calculate Gini Index of each attribute using below equation

$$Gini(Target, Attribute) = \sum_{Branch \in Attribute} p(branch)Gini(branch)$$

$$where Gini(branch) = Gini(a_1, a_2, \dots, a_n) = 1 - \sum_{i \in n} \left[\frac{a_i}{\sum_{i=1}^n a_i} \right]^2$$

3. Choose attribute with lowest Gini as decision node,

Repeat step 2 and 3 until complete tree formed with leaf nodes