

WEEK- 7: Decision trees.

Session No.4

Decision Tree

- Decision tree algorithm falls under the category of supervised learning. They can be used to solve both regression and classification problems.
- Decision tree uses the tree representation to solve the problem.
- A decision tree is a classification and prediction tool having a tree-like structure, where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (terminal node) holds a class label.

Below are some assumptions that we made while using decision tree:

- At the beginning, we consider the whole training set as the root.
- Feature values are preferred to be categorical. If the values are continuous then they are discretized prior to building the model.
- On the basis of attribute values records are distributed recursively.
- We use statistical methods for ordering attributes as root or the internal node.

In Decision Tree the major challenge is to identify the attribute for the root node in each level. This process is known as attribute selection. We have two popular attribute selection measures:

1. Information Gain
2. Gini Index

Attribute selection using Information Gain

Entropy:

Entropy is a measure of the randomness in the information being processed. The higher the entropy, the harder it is to draw any conclusions from that information.

$$H(X) = - \sum_{i=1}^n P(x_i) \log_b P(x_i)$$

Information Gain:

Information gain can be defined as the amount of information gained about a random variable or signal from observing another random variable. It can be considered as the difference between the entropy of parent node and weighted average entropy of child nodes.

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

Alternatively,

$$IG(S, A) = H(S) - \sum_{i=0}^n P(x) * H(x)$$

Building Decision Tree using Information Gain

- Start with all training instances associated with the root node
- Use information gain to choose which attribute to label each node with
- Note: No root-to-leaf path should contain the same discrete attribute twice
- Recursively construct each subtree on the subset of training instances that would be classified down that path in the tree.

Example:

Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No

Find the entropy of the class variable.

$$E(S) = -[(9/14)\log(9/14) + (5/14)\log(5/14)] = 0.94$$

where, S = Playing Football

From the above data for outlook, we can arrive at the following table

		play		total
		yes	no	
Outlook	sunny	3	2	5
	overcast	4	0	4
	rainy	2	3	5
				14

Now we have to calculate average weighted entropy. ie, we have found the total of weights of each feature multiplied by probabilities.

$E(S, \text{outlook})$

$$\begin{aligned}
 &= P(\text{Sunny}) * \text{Entropy}(\text{Sunny}) + P(\text{Overcast}) * \text{Entropy}(\text{Overcast}) + P(\text{Rain}) * \text{Entropy}(\text{Rain}) \\
 &= (5/14) * E(2,3) + (4/14) * E(4,0) + (5/14) * E(3,2) \\
 &= (5/14) [-(2/5) \log(2/5) - (3/5) \log(3/5)] + (4/14)(0) + (5/14) [-(3/5) \log(3/5) - (2/5) \log(2/5)] \\
 &= 0.693
 \end{aligned}$$

The next step is to find the information gain. It is the difference between parent entropy and average weighted entropy we found above.

$$IG(S, \text{outlook}) = 0.94 - 0.693 = 0.247$$

Similarly find Information gain for Temperature, Humidity, and Windy. $IG(S, \text{Temperature}) = 0.940 - 0.911 = 0.029$

$$IG(S, \text{Humidity}) = 0.940 - 0.788 = 0.152$$

$$IG(S, \text{Windy}) = 0.940 - 0.8932 = 0.048$$

Now select the feature having the largest entropy gain. Here it is Outlook. So it forms the first node (root node) of our decision tree.

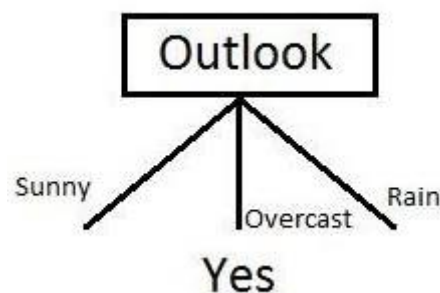
Now our data look as follows

Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes

Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Overcast	Hot	High	Weak	Yes
Overcast	Cool	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes

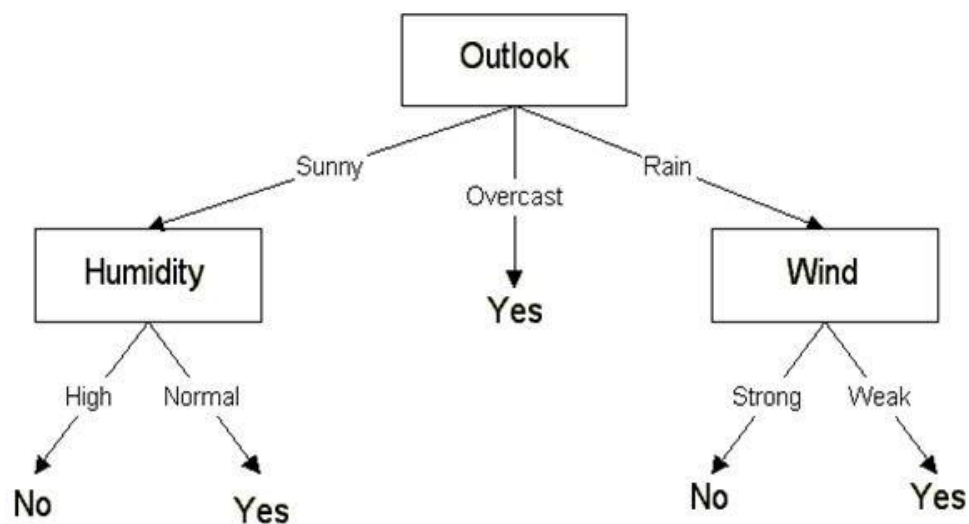
Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Rain	Mild	Normal	Weak	Yes
Rain	Mild	High	Strong	No

Since overcast contains only examples of class 'Yes' we can set it as yes. That means If outlook is overcast football will be played. Now our decision tree looks as follows.



Repeat the above process for each type under Outlook, i.e., Sunny and Rain to find the attribute with highest Information gain.

At the end, our decision tree will look as below:



Overfitting and Decision Trees

Decision Trees are prone to over-fitting. A decision tree will always overfit the training data if we allow it to grow to its max depth. A decision tree is overfit when the tree is trained to fit all samples in the training data set perfectly.

You can tweak some parameters such as min_samples_leaf to minimize default overfitting. The deeper you allow your tree to grow, the more complex the sequence of decision rules becomes. Assigning a **maximum depth** to a tree can simplify it and combat overfitting.

Pruning

Pruning is a technique that is used to reduce overfitting. Pruning also simplifies a decision tree by removing the weakest rules. Pruning is often distinguished into:

- **Pre-pruning** (early stopping) stops the tree before it has completed classifying the training set,
- **Post-pruning** allows the tree to classify the training set perfectly and then prunes the tree.

Pre-pruning

- The pre-pruning technique of Decision Trees is tuning the hyperparameters prior to the training pipeline.
- It involves the heuristic known as 'early stopping' which stops the growth of the decision tree - preventing it from reaching its full depth.
- It stops the tree-building process to avoid producing leaves with small samples.
- During each stage of the splitting of the tree, the cross-validation error will be monitored.
- If the value of the error does not decrease anymore - then we stop the growth of the decision tree.
- The hyperparameters that can be tuned for early stopping and preventing overfitting are: max_depth, min_samples_leaf and min_samples_split

Post-pruning

- Post-pruning does the opposite of pre-pruning and allows the Decision Tree model to grow to its full depth.
- Once the model grows to its full depth, tree branches are removed to prevent the model from overfitting.
- The algorithm will continue to partition data into smaller subsets until the final subsets produced are similar in terms of the outcome variable.
- The final subset of the tree will consist of only a few data points allowing the tree to have learned the data to the T.
- However, when a new data point is introduced that differs from the learned data - it may not get predicted well.
- The hyperparameter that can be tuned for post-pruning and preventing overfitting is ccp_alpha

- ccp stands for Cost Complexity Pruning and can be used as another option to control the size of a tree. A higher value of `ccp_alpha` will lead to an increase in the number of nodes pruned.

