

Introduction to Decision Trees

Friday, August 23, 2019 9:09 PM

Decision Trees -

1. Decision Tree is one of the most commonly used models in data science world
2. It also is a proven management tool used to take decisions in complex situations
3. It can be used for regression and classification, more often used for classification
4. Can be used for binary classification such as whether an applicant for loan is likely to turn into defaulter or not, whether a customer is likely to churn or not
5. It can also be used for multi-class classification for e.g. identifying the character in English alphabet
6. Decision Tree algorithm finds the relation between the target column and the independent variables and expresses it as a tree structure
7. It does so by binary splitting data using functions based on comparison operators on the independent columns

Decision Trees Structure – Training / Building

1. Suppose we are given the data about cars as shown
2. Our objective is to find if any patterns exist that connect the "Horse-Power" and "Weight" to car type (Large or Small)
3. The independent variables are "Horse-Power" and "Weight" while the target column is "Car Type"
4. The target column has binary values (L and S) in equal numbers.

Horse-Power	Weight	Car Type
130	3500	L
90	2000	S
90	1500	S
150	3000	L
270	2500	L
200	2900	L
70	2530	S
215	2000	L
80	2200	S
100	1700	S

Decision Trees Structure – Training / Building

Let us apply a function on the weight column.

Horse-Power	Weight	Car Type
130	3500	L
90	2000	S
90	1500	S
150	3000	L
270	2500	L
200	2900	L
70	2530	S
215	2000	L
80	2200	S
100	1700	S

Horse-Power	Weight	Car Type
80	2200	S
270	2500	L
70	2530	S
200	2900	L
150	3000	L
130	3500	L

Weight > 2000.

Horse-Power	Weight	Car Type
90	1500	S
100	1700	S
90	2000	S
215	2000	L

1. This smaller node on top has "L" in majority in the target column hence gets label "L"
2. The smaller node on the bottom has "S" in majority in the target column hence gets label "S"
3. The homogeneity of the target column in both the smaller nodes has increased compared to parent
4. But both the smaller nodes still have a mix of values in the target column
5. Let us split the data further using Horse_Power

Decision Trees Structure – Training / Building

Class = L (Majority)

Horse-Power	Weight	Car Type
80	2200	S
270	2500	L
70	2530	S
200	2900	L
150	3000	L
130	3500	L

Let us apply a function on the Horsepower column

Horse-Power	Weight	Car Type
90	1500	S
100	1700	S
90	2000	S
215	2000	L

Class = S (Majority)

Horse-Power	Weight	Car Type
70	2530	S
80	2200	S
80	2200	S

HorsePower > 100.

Horse-Power	Weight	Car Type
150	3000	L
200	2900	L
270	2500	L

Horse-Power	Weight	Car Type
90	1500	S
90	2000	S
100	1700	S

HorsePower > 200.

Horse-Power	Weight	Car Type
215	2000	L

1. The smaller node on the top now is perfectly homogenous in target column and belongs to class "S"
2. The second node similarly belongs to "L"
3. The third node belongs to "S"
4. The fourth node belongs to "L"
5. There is no further need to split the data as it is perfectly homogenous!

Decision Trees Structure – Training / Building

Let us apply a function on the weight column.

Horse-Power	Weight	Car Type
130	3500	L
90	2000	S
90	1500	S
150	3000	L
270	2500	L
200	2900	L
70	2530	S
215	2000	L
90	2200	S
100	1700	S



The tree thus has given us the following relation between "Weight", "Horse-Power" and "Car Type" as

If wt > 2000
if hp > 100
class = "L"
else
class = "S"

If wt <= 2000
if hp > 200
class = "L"
else
class = "S"

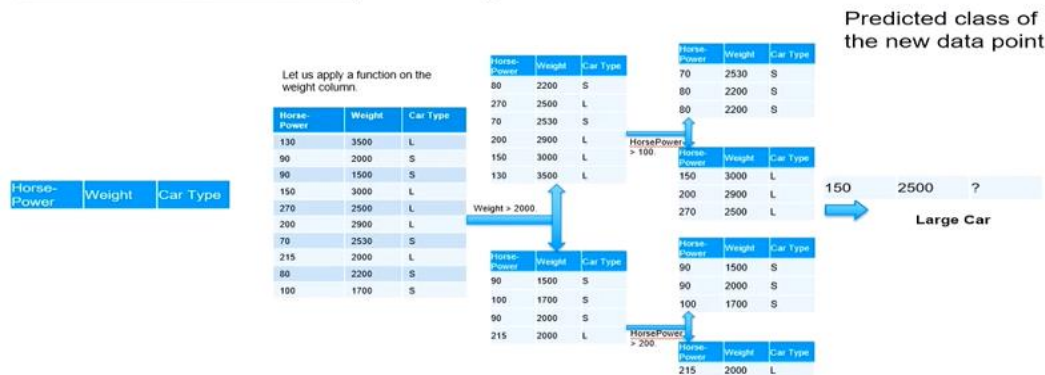
Note: The CART algorithm employed by scikitlearn creates only binary tree i.e. each node is split into two subnodes

Decision Trees Structure - Predicting

Predicted class of the new data point



Decision Trees Structure - Predicting



Decision Trees Structure – Training Errors

1. Suppose we come across a combination of "Weight" and "Horse-Power" for any of the classes which was not available in the training data on which the decision tree was built. For e.g. a Small Car with HorsePower of 250 and Weight is 2000. The Decision Tree will classify it as Large Car.
2. Such classification errors can occur both during the training and testing. They are called training errors and testing errors. This is true for any algorithm
3. The decision tree algorithm by default will try to build a tree where the smallest child nodes are perfectly homogenous in the target columns
4. To achieve perfect homogeneity in the target column, the algorithm may build a large tree where each leaf has only 1 record!!! Such models are overfit models. They give zero errors on training but perform poorly on test data

Decision Trees – Posterior Probability

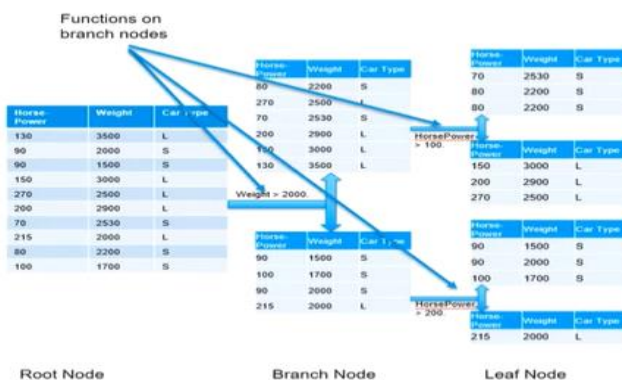
5. Sometimes when the algorithm runs out of independent attributes to use to break a node into smaller nodes or it is forced to stop, we may find nodes where the target column is not homogenous

Small Cars

Horse-Power	Weight	Car Type
90	1500	S
100	1700	S
90	2000	S
215	2000	L

6. In such case, the label assigned to the node is based on majority class and the ratio of the classes indicates the posterior probability of the two classes at that node. $P(S) = \frac{3}{4}$ and $P(L) = \frac{1}{4}$

Decision Trees – Structure & Node types



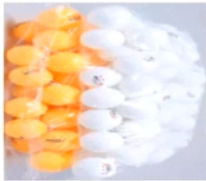
7. Decision trees consist of nodes and branches. Nodes represent a decision function while branch represents the result of the function. Thus it is a flow chart for deciding how to classify a new observation:
8. The nodes are of three types, Root Node (representing the original data), Branch Node (representing a function), Leaf Node (which holds the result of all the previous functions that connect to it)

How does decision tree algorithm select the column and the breakpoint in the column to create the tree?

Decision Trees – Learning Process / Loss function & impurity

1. The decision tree algorithm learns (i.e. creates the decision tree from the data set) through optimization of a loss function
2. The loss function represents the reduction in impurity of the target column i.e. increase in the homogeneity of the target column at every split of the given data
3. To understand the loss function we need to understand how it computes the impurity / purity of the target column. There are two measures of impurity viz Entropy and Gini...

Decision Trees – Learning Process / Loss function & impurity



Decision Trees – Learning Process / Loss function & impurity



1. There is a bag of 50 balls of orange and white respectively
2. You have to pull out one ball from the bag with closed eyes. If the ball is -
 - a. orange, explain the linear regression model
 - b. white, nothing to be done
3. This state where you have to decide and your decision can result in multiple outcomes with equal probability is said to be state of maximum uncertainty



4. If you have a bag full of balls of only one colour, then there is no uncertainty. You know what is going to happen. Uncertainty is zero.
5. Thus, the more the homogeneity, lesser the uncertainty and vice versa
6. Uncertainty is expressed as entropy or Gini index

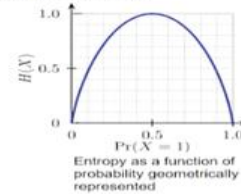
Entropy ... a measure of uncertainty

- a. Given that there are two possible outcomes (orange or white) for a given action (picking the ball)
- b. Given that we know the probability of getting each outcome ($P(\text{Orange}) = .5$ and $P(\text{White}) = 1 - .5 = .5$)
- c. We also know that when only one outcome (White or Orange) is possible, there is no uncertainty. $P(\text{White} = 1)$ and $P(\text{Orange} = 1 - 1 = 0)$
- d. We can express the relation between probability and impurity of target column in a mathematical form

Decision Trees – Shannon's Entropy

- Imagine a bag contains 6 Orange and 4 White balls. Let the two classes Orange -> class 0 and White -> class 1. The probability ranges from 0 to 1
- The impurity of the bag can be represented as a log to base 2 of probability of a class (pi). Impurity ranges from 0 to 1 (for binary classification). The relation between probability and impurity can be expressed as -

$$H(X) = - \sum_{i=0}^{N-1} p_i \log_2 p_i$$



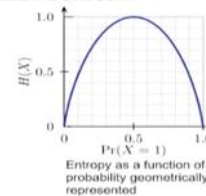
- H(X), entropy of the bag (X) will be calculated as per the formula
 - The index i refers to the number of classes possible
 - $H(X) = -(0.6 * \log_2(0.6)) - (0.4 * \log_2(0.4)) = 0.9709506$

- Suppose we remove all Orange balls from the bag and then entropy will be
 - $H(X) = -1.0 * \log_2(1.0) - 0.0 * \log_2(0) = 0$ ## Entropy is 0! i.e. Information is 100%

Decision Trees – Shannon's Entropy

- Imagine a bag contains 6 Orange and 4 White balls. Let the two classes Orange -> class 0 and White -> class 1. The probability ranges from 0 to 1
- The impurity of the bag can be represented as a log to base 2 of probability of a class (pi). Impurity ranges from 0 to 1 (for binary classification). The relation between probability and impurity can be expressed as -

$$H(X) = - \sum_{i=0}^{N-1} p_i \log_2 p_i$$



- H(X), entropy of the bag (X) will be calculated as per the formula
 - The index i refers to the number of classes possible
 - $H(X) = -(0.6 * \log_2(0.6)) - (0.4 * \log_2(0.4)) = 0.9709506$

- Suppose we remove all Orange balls from the bag and then entropy will be
 - $H(X) = -1.0 * \log_2(1.0) - 0.0 * \log_2(0) = 0$ ## Entropy is 0! i.e. Information is 100%

Decision Trees – Entropy & Objective

Suppose we wish to find if there was any influence of shipping mode, order priority on customer location. Customer location is target column and like the bag of coloured balls

Row #	Order ID	Order Date	Order Priority	Order Quantity	Sales	Discount	Ship Mode	Profit	Unit Price	Shipping Cost	Customer Name	Province
1	1	13-10-2010	Low	6	261.54	0.04	Regular Air	-213.25	36.94	36	Mohamed Mactay	Namur
49	293	01-10-2012	High	49	10123.02	0.07	Delivery Truck	467.81	208.16	68.02	Berry French	Namur
50	293	01-10-2012	High	27	244.57	0.01	Regular Air	46.71	8.69	2.99	Berry French	Namur
80	403	10-01-2011	High	30	4905.7595	0.08	Regular Air	1198.57	195.99	3.99	Clay Piccardi	Namur
85	515	20-06-2010	Not Specified	19	394.27	0.08	Regular Air	30.94	21.70	5.94	Carlos Soltero	Namur
86	515	20-06-2010	Not Specified	21	146.69	0.05	Regular Air	4.43	6.64	4.95	Carlos Soltero	Namur
97	613	17-06-2011	High	12	93.54	0.03	Regular Air	-54.04	7.3	7.72	Carl Jackson	Namur
98	613	17-06-2011	High	22	905.08	0.09	Regular Air	127.70	42.76	6.22	Carl Jackson	Namur
103	643	24-03-2011	High	21	2781.82	0.07	Express Air	495.26	138.14	35	Monica Faderle	Namur
107	678	26-02-2010	Low	44	228.41	0.07	Regular Air	-226.36	4.98	6.33	Conorhy Badders	Namur
127	807	23-11-2010	Medium	45	196.05	0.01	Regular Air	-168.85	4.20	6.18	Nicola Schneider	Namur
128	807	23-11-2010	Medium	32	124.56	0.04	Regular Air	-14.33	3.95	2	Nicola Schneider	Namur
134	868	08-06-2012	Not Specified	32	716.84	0	Regular Air	134.72	21.70	5.94	Carlos Daly	Namur
135	868	08-06-2012	Not Specified	31	1474.33	0.04	Regular Air	114.46	47.80	3.61	Carlos Daly	Namur
149	933	04-09-2010	Regular Air	15	88.61	0.02	Regular Air	-4.72	5.20	2.99	Claudio Miner	Namur
160	995	30-05-2011	Medium	46	1015.49	0.03	Regular Air	792.91	39.89	3.04	Nicola Schneider	Namur
161	998	25-11-2009	Not Specified	16	240.26	0.07	Regular Air	93.80	15.74	1.39	Allen Posenblatt	Namur
175	1154	14-02-2012	Critical	44	4482.23	0.04	Delivery Truck	440.72	188.98	26.22	Sylvia Foulton	Namur
176	1154	14-02-2012	Critical	11	663.784	0.25	Regular Air	-481.84	71.27	69	Sylvia Foulton	Namur
203	1344	15-04-2012	Low	15	834.904	0.06	Regular Air	-11.68	65.99	5.26	Jim Redford	Namur
204	1344	15-04-2012	Low	18	2488.9205	0.01	Regular Air	313.58	155.99	8.99	Jim Redford	Namur
213	1402	12-03-2010	Not Specified	13	59.03	0.1	Express Air	26.92	3.69	0.5	Carlos Soltero	Namur
214	1402	12-03-2010	Not Specified	21	97.40	0.05	Regular Air	-7.77	4.71	0.7	Carlos Soltero	Namur
229	1539	09-03-2011	Low	33	511.03	0.1	Regular Air	-172.88	35.99	13.18	Carl Ludwig	Namur
230	1539	09-03-2011	Low	38	194.99	0.05	Regular Air	-144.95	4.89	4.93	Carl Ludwig	Namur
231	1540	04-09-2012	High	30	80.9	0.1	Regular Air	5.76	2.88	0.7	Con Miller	Namur
249	1702	06-05-2011	High	23	67.24	0.08	Regular Air	4.93	2.84	0.93	Anna Cyprus	Namur

Sales Data

Decision Trees – Entropy & Objective

Suppose we wish to find if there was any influence of shipping mode, order priority on customer location. Customer location is target column and like the bag of coloured balls

Row #	Order #	Order Date	Order Priority	Order Quantity	Sales	Discount	Ship Mode	Profit	Unit Price	Shipping Cost	Customer Name	Province
733	5678	10-11-2019	High	7	754.32	0	Regular Air	-123.57	99.99	13.99	Adam Shams	Alberta
8867	57216	29-07-2019	High	46	109.76	0.07	Regular Air	4.46	2.2	1.2	Ruben Darr	Alberta
8017	57201	29-04-2019	High	26	561.68	0.07	Regular Air	-111.33	22.38	15.1	Andy Yellou	Alberta
886	5767	29-04-2019	High	31	2780.54	0.08	Regular Air	786.10	82.99	5.5	Jessica Myrick	Alberta
1	1	13-10-2019	Low	6	261.54	0.04	Regular Air	-113.25	38.94	35	Muhammad Moshay	Manitoba
167	678	28-02-2019	Low	44	228.47	0.07	Regular Air	-28.36	4.98	8.13	Doreilly Bashley	Manitoba
231	1344	16-04-2012	Low	15	834.804	0.06	Regular Air	-11.69	88.99	5.26	Jim Radford	Manitoba
234	1344	16-04-2012	Low	18	2480.9205	0.07	Regular Air	313.59	105.99	9.99	Jim Radford	Manitoba
229	1538	09-03-2011	Low	32	511.83	0.1	Regular Air	-172.86	15.99	13.19	Carl Ludwig	Manitoba
236	1538	09-03-2011	Low	36	164.99	0.08	Regular Air	-144.95	4.99	4.93	Carl Ludwig	Manitoba
236	1762	08-11-2019	Low	29	370.46	0.04	Regular Air	-6.45	13.46	4.93	Carlos Sotero	Manitoba
381	3207	23-09-2019	Low	27	1078.49	0.08	Regular Air	252.96	46.96	1.89	Jack Costa	Manitoba
755	5469	09-07-2012	Low	11	46.91	0.07	Regular Air	-7.04	3.98	2.97	Clon Jones	Manitoba
921	9205	12-10-2017	Low	44	3922.42	0.08	Regular Air	-954.96	92.23	39.97	Jay Gell	Manitoba



Row #	Order #	Order Date	Order Priority	Order Quantity	Sales	Discount	Ship Mode	Profit	Unit Price	Shipping Cost	Customer Name	Province
4885	4039	13-09-2019	Not Specified	4	394	0.01	Express Air	286.31	92.99	10.94	Victoria Brennan	Alberta
7919	56728	31-10-2017	Low	21	80.88	0.07	Express Air	36.81	4.71	8.91	Ruben Darr	Alberta
7989	15822	23-04-2012	Medium	11	821.496	0	Express Air	128.96	46.99	1.25	Shu Tan	Alberta
7688	54612	03-05-2019	High	11	5432	0.06	Express Air	1439	4.76	6.88	Tanya Turrel	Alberta
7758	95424	00-05-2019	Medium	6	25.84	0.04	Express Air	-41.13	5.76	4.93	Victoria Brennan	Alberta
7688	56727	11-01-2019	Medium	47	1076.75	0.08	Express Air	357.23	29.99	9.99	Victoria Brennan	Alberta
885	5767	29-04-2019	High	36	163.94	0.03	Express Air	-65.96	4.71	5.94	Jessica Myrick	Alberta

Decision Trees – Entropy & Objective

Suppose we wish to find if there was any influence of shipping mode, order priority on customer location. Customer location is target column and like the bag of coloured balls

Row #	Order #	Order Date	Order Priority	Order Quantity	Sales	Discount	Ship Mode	Profit	Unit Price	Shipping Cost	Customer Name	Province
733	5678	10-11-2019	High	7	754.32	0	Regular Air	-123.57	99.99	13.99	Adam Shams	Alberta
8867	57216	29-07-2019	High	46	109.76	0.07	Regular Air	4.46	2.2	1.2	Ruben Darr	Alberta
8017	57201	29-04-2019	High	26	561.68	0.07	Regular Air	-111.33	22.38	15.1	Andy Yellou	Alberta
886	5767	29-04-2019	High	31	2780.54	0.08	Regular Air	786.10	82.99	5.5	Jessica Myrick	Alberta



Row #	Order #	Order Date	Order Priority	Order Quantity	Sales	Discount	Ship Mode	Profit	Unit Price	Shipping Cost	Customer Name	Province
1	1	13-10-2019	Low	6	261.54	0.04	Regular Air	-113.25	38.94	35	Muhammad Moshay	Manitoba
167	678	28-02-2019	Low	44	228.47	0.07	Regular Air	-28.36	4.98	8.13	Doreilly Bashley	Manitoba
231	1344	16-04-2012	Low	15	834.804	0.06	Regular Air	-11.69	88.99	5.26	Jim Radford	Manitoba
234	1344	16-04-2012	Low	18	2480.9205	0.07	Regular Air	313.59	105.99	9.99	Jim Radford	Manitoba
229	1538	09-03-2011	Low	32	511.83	0.1	Regular Air	-172.86	15.99	13.19	Carl Ludwig	Manitoba
236	1538	09-03-2011	Low	36	164.99	0.08	Regular Air	-144.95	4.99	4.93	Carl Ludwig	Manitoba
236	1762	08-11-2019	Low	29	370.46	0.04	Regular Air	-6.45	13.46	4.93	Carlos Sotero	Manitoba
381	3207	23-09-2019	Low	27	1078.49	0.08	Regular Air	252.96	46.96	1.89	Jack Costa	Manitoba
755	5469	09-07-2012	Low	11	46.91	0.07	Regular Air	-7.04	3.98	2.97	Clon Jones	Manitoba

Decision Trees – Entropy & Objective

Suppose we wish to find if there was any influence of shipping mode, order priority on customer location. Customer location is target column and like the bag of coloured balls

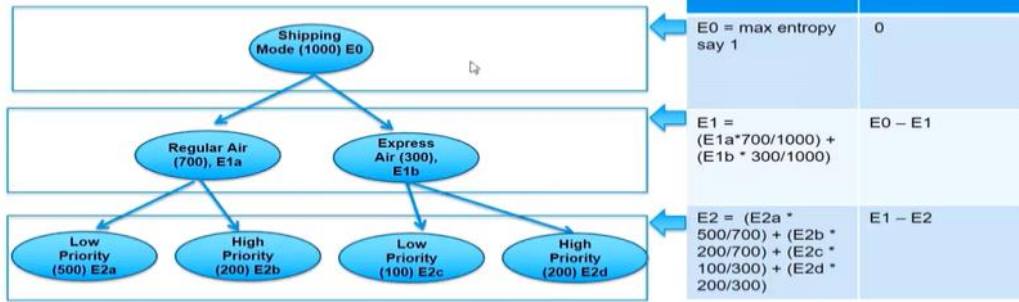
Row #	Order #	Order Date	Order Priority	Order Quantity	Sales	Discount	Ship Mode	Profit	Unit Price	Shipping Cost	Customer Name	Province
733	5678	10-11-2019	High	7	754.32	0	Regular Air	-123.57	99.99	13.99	Adam Shams	Alberta
8867	57216	29-07-2019	High	46	109.76	0.07	Regular Air	4.46	2.2	1.2	Ruben Darr	Alberta
8017	57201	29-04-2019	High	26	561.68	0.07	Regular Air	-111.33	22.38	15.1	Andy Yellou	Alberta
886	5767	29-04-2019	High	31	2780.54	0.08	Regular Air	786.10	82.99	5.5	Jessica Myrick	Alberta



Row #	Order #	Order Date	Order Priority	Order Quantity	Sales	Discount	Ship Mode	Profit	Unit Price	Shipping Cost	Customer Name	Province
733	5678	10-11-2019	High	7	754.32	0	Regular Air	-123.57	99.99	13.99	Adam Shams	Alberta
8867	57216	29-07-2019	High	46	109.76	0.07	Regular Air	4.46	2.2	1.2	Ruben Darr	Alberta
8017	57201	29-04-2019	High	26	561.68	0.07	Regular Air	-111.33	22.38	15.1	Andy Yellou	Alberta
886	5767	29-04-2019	High	31	2780.54	0.08	Regular Air	786.10	82.99	5.5	Jessica Myrick	Alberta

When sub branches are created, the total entropy of the sub branches should be less than the entropy of the parent node. More the drop in entropy, more the information gained

Decision Trees – Entropy & Objective



Tree will stop growing when stop criterion for the splitting is reached which could be -

- Tree has reached certain pre-fixed depth (longest path from root node to leaf node)
- Tree has achieved maximum number of nodes (tree size)
- Exhausted all attributes to split
- Leaf node on split will have less than predefined number of data points

Decision Trees – Loss functions

Common measures of purity

- Gini index – is calculated by subtracting the sum of the squared probabilities of each class from one

- Uses squared proportion of classes
- Perfectly classified, Gini Index would be zero
- Evenly distributed would be $1 - (1/\# \text{ Classes})$
- You want a variable split that has a low Gini Index
- Used in CART algorithm

$$Gini = 1 - \sum_{i=1}^C (p_i)^2$$

- Entropy –

- Favors splits with small counts but many unique value
- Weights probability of class by $\log(\text{base}=2)$ of the class probability
- A smaller value of Entropy is better. That makes the difference between the parent node's entropy larger
- Information Gain is the Entropy of the parent node minus the entropy of the child nodes

$$Entropy = \sum_{i=1}^C -p_i * \log_2(p_i)$$

Decision Trees - Information Gain using Entropy

$$H(X) = - \sum_{i=0}^{N-1} p_i \log_2 p_i$$

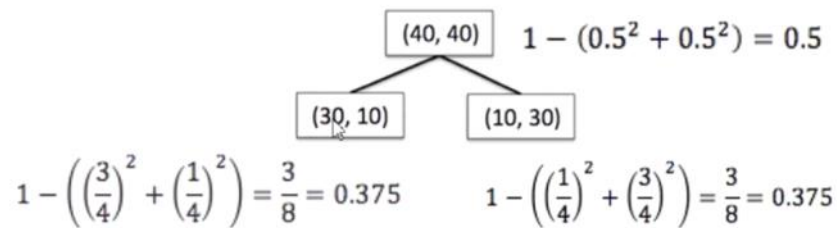
$$\begin{array}{c}
 (40, 40) \\
 \swarrow \quad \searrow \\
 (30, 10) \quad (10, 30)
 \end{array}
 \quad
 I_H(D_p) = -(0.5 \log_2(0.5) + 0.5 \log_2(0.5)) = 1$$

$$\begin{array}{c}
 -\left(\frac{3}{4} \log_2\left(\frac{3}{4}\right) + \frac{1}{4} \log_2\left(\frac{1}{4}\right)\right) = 0.81 \\
 -\left(\frac{1}{4} \log_2\left(\frac{1}{4}\right) + \frac{3}{4} \log_2\left(\frac{3}{4}\right)\right) = 0.81
 \end{array}$$

$$\text{Information Gain} = \text{reduction in entropy} = 1 - \frac{4}{8} 0.81 - \frac{4}{8} 0.81 = 0.19$$

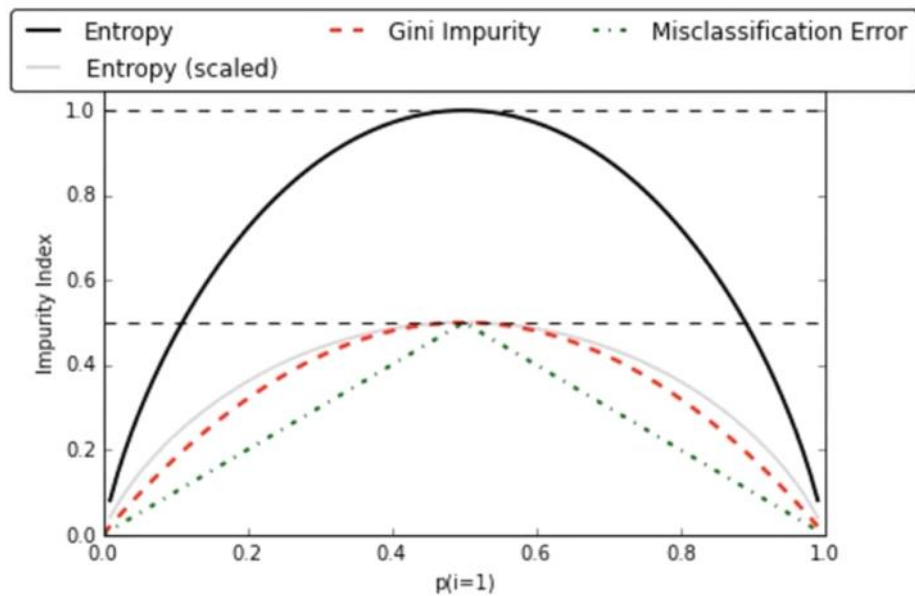
Decision Trees - Information Gain using Gini index

$$Gini = 1 - \sum_{i=1}^C (p_i)^2$$



$$\text{Information Gain} = \text{reduction in Gini index} = 0.5 - \frac{4}{8} \cdot 0.375 - \frac{4}{8} \cdot 0.375 = 0.125$$

Decision Trees – Gini , Entropy , Misclassification Error



Note: Misclassification Error is not used in Decision Trees

Decision Trees -

Advantages -

1. Simple , Fast in processing and effective
2. Can work with missing data
3. Handles numeric and categorical variables
4. Interpretation of results is easier when represented as rules

Dis-advantages -

1. Often biased towards features have large number of levels
2. May not be optimum as modelling some relations on axis parallel basis is not optimal
3. Small changes in training data can result in large changes to the logic
4. Decision trees tend to become overfit by default creating large complex trees
5. Large trees can be difficult to interpret

Decision Trees - Algorithms

1. ID3 (Iterative Dicotomizer 3) – developed by Ross Quinlan. Creates a multi branch tree at each node using greedy algorithm. Trees grow to maximum size before pruning
2. C4.5 succeeded ID3 by overcoming limitation of features required to be categorical. It dynamically defines discrete attribute for numerical attributes. It converts the trained trees into a set of if-then rules. Accuracy of each rule is evaluated to determine the order in which they should be applied
3. C5.0 is Quinlan's latest version and it uses less memory and builds smaller rulesets than C4.5 while being more accurate
4. CART (Classification & Regression Trees) is similar to C4.5 but it supports numerical target variables and does not compute rule sets. Creates binary tree. Scikit uses CART

Help

Friday, August 23, 2019 9:21 PM

<https://scikit-learn.org/stable/modules/tree.html>