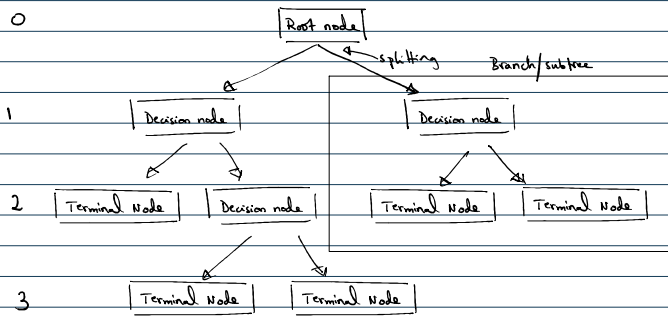


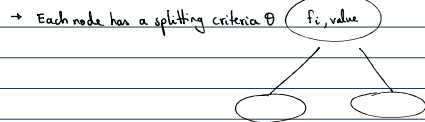
→ Non-Parametric model
→ gaint Nested if-else structure → DT

Terminology



CART (Classification and Regression Trees) (Binary trees)

→ Algorithm to grow decision tree



→ Which feature & value to perform splitting?

Algorithm

Given training vectors $x_1 \in \mathbb{R}^P, \dots, x_n \in \mathbb{R}^P$ and a label vector $y \in \mathbb{R}^P$, a decision tree recursively partitions the feature space such that the samples with the same labels or similar target values are grouped together.

Let the data at node m be represented by Q_m with n_m samples. For each candidate split $\theta = (j, t_m)$ consisting of a feature j and threshold t_m , partition the data into $Q_m^{left}(\theta)$ and $Q_m^{right}(\theta)$ subsets

$$Q_m^{left}(\theta) = \{(x, y) | x_j \leq t_m\}$$
$$Q_m^{right}(\theta) = Q_m \setminus Q_m^{left}(\theta)$$

The quality of a candidate split of node m is then computed using an impurity function or loss function $H()$, the choice of which depends on the task being solved (classification or regression)

$$G(Q_m, \theta) = \frac{n_m^{left}}{n_m} H(Q_m^{left}(\theta)) + \frac{n_m^{right}}{n_m} H(Q_m^{right}(\theta))$$

Select the parameters that minimises the impurity

$$\theta^* = \underset{\theta}{\operatorname{argmin}} G(Q_m, \theta)$$

Recurse for subsets $Q_m^{left}(\theta^*)$ and $Q_m^{right}(\theta^*)$ until the maximum allowable depth is reached, $n_{max} < \min_{\text{samples}} \text{ or } n_{max} = 1$.

- Splitting Criteria:** At each step of building the tree, the algorithm looks for the best way to split the data. It does this by choosing a feature and a value to split on. The goal is to find the split that maximizes the homogeneity of the target variable within each subset created by the split.
- Impurity Measures:** Sklearn uses impurity measures like Gini impurity or entropy to quantify the homogeneity within a set of data points. Gini impurity measures the probability of incorrectly classifying a randomly chosen element if it was randomly labeled according to the distribution of labels in the subset. Entropy, on the other hand, measures the average amount of information needed to classify a data point.
- Optimal Split Selection:** The algorithm evaluates different splits by calculating the impurity of the subsets created by each split. It then chooses the split that minimizes the impurity or maximizes the information gain (the difference between the impurity before and after the split).
- Recursive Partitioning:** This process of selecting the best split and partitioning the data is repeated recursively for each subset until certain stopping criteria are met, such as reaching a maximum tree depth, having too few data points in a node, or no further improvement in impurity reduction.
- Prediction:** Once the tree is built, to make predictions for new data points, the algorithm navigates through the tree based on the feature values of the data point until it reaches a leaf node. The majority class or the average value of the target variable in that leaf node is then assigned as the predicted value for the new data point.

Recursive Partitioning Algorithm:-

Suppose we have a response variable Y and a set of P predictor variables X_j for $j = 1, \dots, P$. For a partition A of records, recursive partitioning will find the best way to partition A into two subpartitions:

- For each predictor variable X_j :
 - For each value s of X_j :
 - Split the records in A with X_j values $< s$ as one partition, and the remaining records where $X_j \geq s$ as another partition.
 - Measure the homogeneity of class within each subpartition of A .
 - Select the value of s that produces maximum within-partition homogeneity of class.
- Select the variable X_j and the split value s_j that produces maximum within-partition homogeneity of class.

Now comes the recursive part.

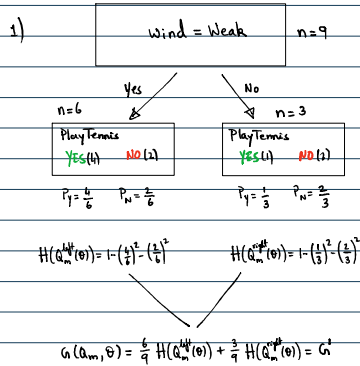
- Initialize A with the entire data set.
- Apply the partitioning algorithm to split A into two subpartitions, A_1 and A_2 .
- Repeat step 2 on subpartitions A_1 and A_2 .
- The algorithm terminates when no further partition can be made that sufficiently improves the homogeneity of the partitions.

Classification problem

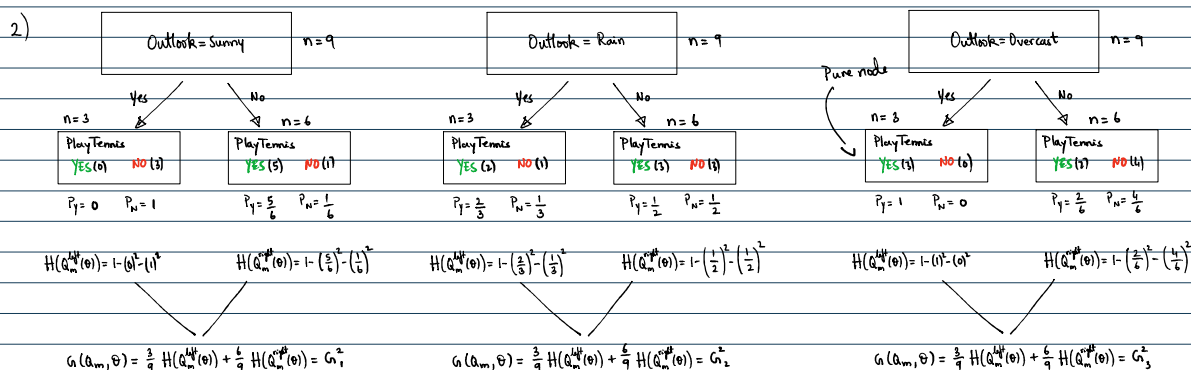
Day	Outlook	Multi class		Numerical		Binary class
		Temperature	Humidity	Wind	PlayTennis	
1	Sunny	85	85	Weak	No	
2	Sunny	80	90	Strong	No	
3	Overcast	83	78	Weak	Yes	
4	Rain	70	96	Weak	Yes	
5	Rain	68	80	Weak	Yes	
6	Rain	65	70	Strong	No	
7	Overcast	64	65	Strong	Yes	
8	Sunny	72	95	Weak	No	
9	Sunny	79	70	Weak	Yes	
10	Rain	75	80	Weak	Yes	
11	Sunny	75	70	Strong	Yes	
12	Overcast	79	90	Strong	Yes	
13	Overcast	83	75	Weak	Yes	
14	Rain	77	80	Strong	No	

$$\text{gini impurity } H() = 1 - \sum_{k=1}^K p_k^2$$

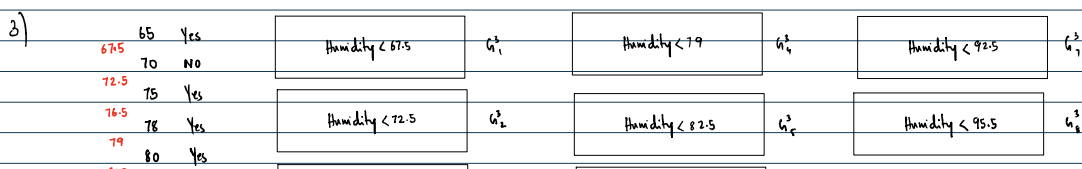
p_k = prob of each class



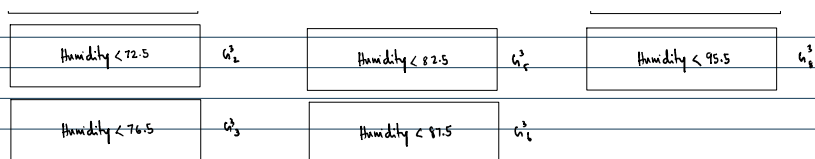
G



$$\min(G_1^2, G_2^2, G_3^2)$$



72.5 75 Yes
 76.5 78 Yes
 79 80 Yes
 82.5 85 No
 87.5 90 No
 92.5 95 No
 95.5 98 Yes



$$\min(G_1, G_2, G_3, G_4, G_5, G_6, G_7, G_8)$$

$$\text{Root node decision} = \operatorname{argmin} \left(\boxed{G_1} \mid \boxed{\min(G_1^1, G_2^1, G_3^1)} \mid \boxed{\min(G_1^3, G_2^3, G_3^3, G_4^3, G_5^3, G_6^3, G_7^3, G_8^3)} \right)$$

4) Repeat until we have pure nodes

Gini Impurity

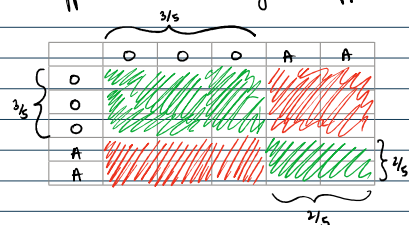
$$H() = 1 - \sum_{i=1}^k p_k^2$$

Gini impurity measures the probability of incorrectly classifying a randomly chosen element if it was randomly labelled according to the distribution of labels in the subset.

Yes, that's a good way to think about it. Gini impurity measures the degree of impurity or uncertainty in a dataset. In a classification problem, if you were to randomly guess the class label for a data point based on the class distribution in the dataset, the Gini impurity represents the probability of guessing incorrectly.

So, when the Gini impurity is low (close to 0), it indicates that the dataset is predominantly pure, and therefore the probability of guessing wrong is low. Conversely, when the Gini impurity is high (close to 1), it suggests that the dataset is more mixed or impure, and the probability of guessing wrong is higher.

→ Suppose we have 3 oranges & 2 apples



$$\text{probability of guessing wrong} = H() = \text{Gini Impurity} = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2$$

↑ Variety ↑ Impurity

$$\text{Gini Impurity Range: } \left[0, 1 - \frac{1}{n}\right]$$

of classes

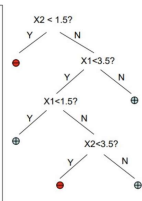
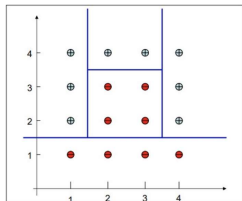
$$0 \leq H() < 1 \rightarrow \text{Categories} \rightarrow \infty$$

→ At each split we aim to attain max purity i.e. min impurity

Geometric perspective

Decision Tree Decision Boundaries

- Decision Trees divide the input space into axis-parallel rectangles and label each rectangle with one of the K classes



Notes

The default values for the parameters controlling the size of the trees (e.g. "max_depth", "min_samples_leaf", etc.) lead to fully grown and unpruned trees which can potentially be very large on some data sets. To reduce memory consumption, the complexity and size of the trees should be controlled by setting those parameter values.

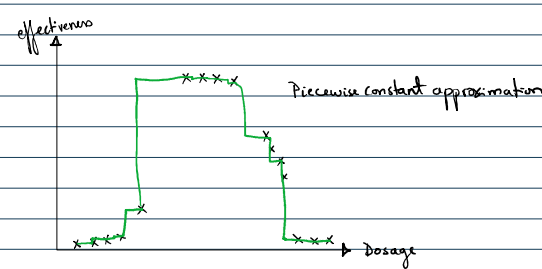
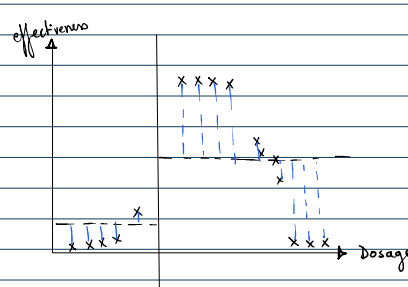
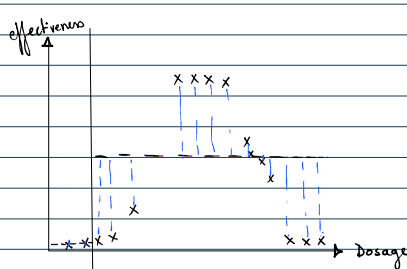
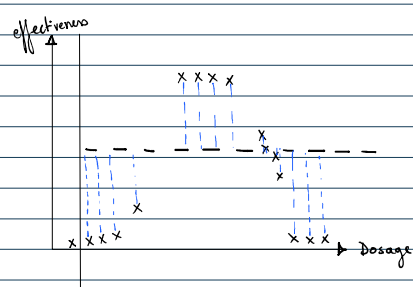
The :meth:'predict' method operates using the :func:'numpy.argmax' function on the outputs of :meth:'predict_proba'. This means that in case the highest predicted probabilities are tied, the classifier will predict the tied class with the lowest index in :term:'classes_'.

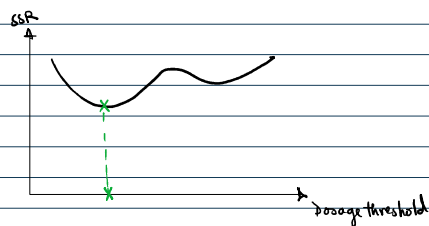
Regression

→ Steps are similar to classification problem, but instead of gini impurity we use Mean Squared Error.

$$\text{mse} = H() = \frac{1}{n} \sum (x_i - \bar{x})^2 \rightarrow \text{Think of this as variance}$$

Geometric perspective





→ Choose threshold value which results in smallest sum of squares residuals.

Pros & Cons

- Simple to understand & interpret. Trees can be visualized
- Requires little data preparation
- The cost of using the tree is logarithmic in the number of data points used to train the tree.
- Able to handle both numerical & categorical data.
- Can work on non-linear datasets
- Can give you feature importance
- DT learners can create over-complex trees that don't generalize the data well. (overfitting)
- DT can be unstable b/c small variations in data might result in a completely different tree being generated.
- predictions of DT's are neither smooth nor continuous, but piecewise constant approximations. ∴ they are not good at extrapolation. ★
- low bias high variance models