

# Decision Trees and Random Forests\*

February 20, 2025

## 1 Introduction

In the previous lectures we dealt with the design of linear classifiers described by linear discriminant functions  $g(\mathbf{x})$ ,

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0.$$

Provided that the dataset is linearly separable, we saw that algorithms such as standard linear classifier and SVM can derive the weights  $\mathbf{w}$  and bias  $b$  which optimally fit the input data samples. When classes slightly overlap, soft margin SVM is exploited. From this lecture, we will deal with problems that are not linearly separable and for which the design of a linear classifier, even in an optimal way, does not lead to satisfactory performance. The design of nonlinear classifiers emerges now as an inescapable necessity.

To seek non-linearly separable problems one does not need to go into complicated situations. The well-known XOR Boolean function is a typical example of such a problem (see Figure 1). Depending on the values of the input binary data  $\mathbf{x} = [x_1, x_2]$ , the output is either 0 or 1, and  $\mathbf{x}$  is classified into one of the two classes  $+1$  or  $-1$ . The corresponding truth table for the XOR operation is shown in Table 1, where a *true* value is expected if the two inputs are not equal and a *false* value if they are equal.

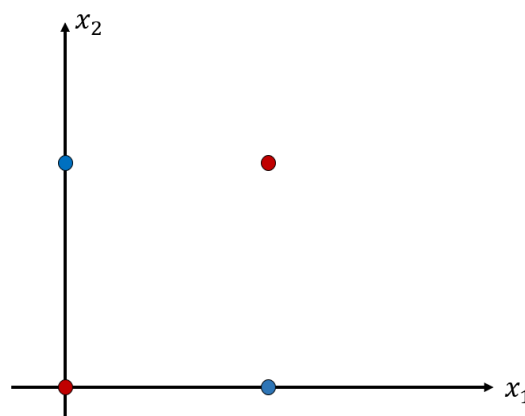


Figure 1: The XOR problem. The blue dots denote *true* and red dots denote *false*.

---

### \*References

- Christopher Bishop. Pattern Recognition and Machine Learning. 2006
- Sergios Theodoridis, Konstantinos Koutroumbas. Pattern Recognition. 2009
- Understanding Random Forests. March 7, 2022

$x_1$	0	0	1	1
$x_2$	0	1	0	1
$y$	+1	-1	-1	+1

Table 1: Truth labels for the XOR problem.

It is apparent from the XOR problem that no single straight line exists that can separate the two classes, as the decision boundary between the two classes is naturally not linear. Non-linear classifiers are specifically designed to cope with such problems. Common non-linear classifiers include decision trees, random forests, Multi-Layer Perceptron (MLP) algorithm, neural networks, and deep learning. In this lecture, we mainly discuss decision trees and random forests.

## 2 Decision trees

Decision trees are known as multistage decision systems in which classes are sequentially rejected until a finally accepted class is reached. To this end, the feature space is split into unique regions, corresponding to the classes, in a sequential manner. Upon the arrival of a feature vector, the searching of the region to which the feature vector will be assigned is achieved via a sequence of decisions along a path of nodes of an appropriately constructed tree. Such schemes offer advantages when a large number of classes are involved. The most popular decision trees are those that split the space into hyperrectangles with sides parallel to the axes. The sequence of decisions is applied to individual features, and the questions to be answered are of the form

$$\text{is feature } x_i \leq a \text{ ?}$$

where  $a$  is a threshold value. Such trees are known as Ordinary Binary Classification Trees (OBCTs). Other types of trees are also possible that split the space into convex polyhedral cells or pieces of spheres. Figure 2 gives an example of the OBCT.

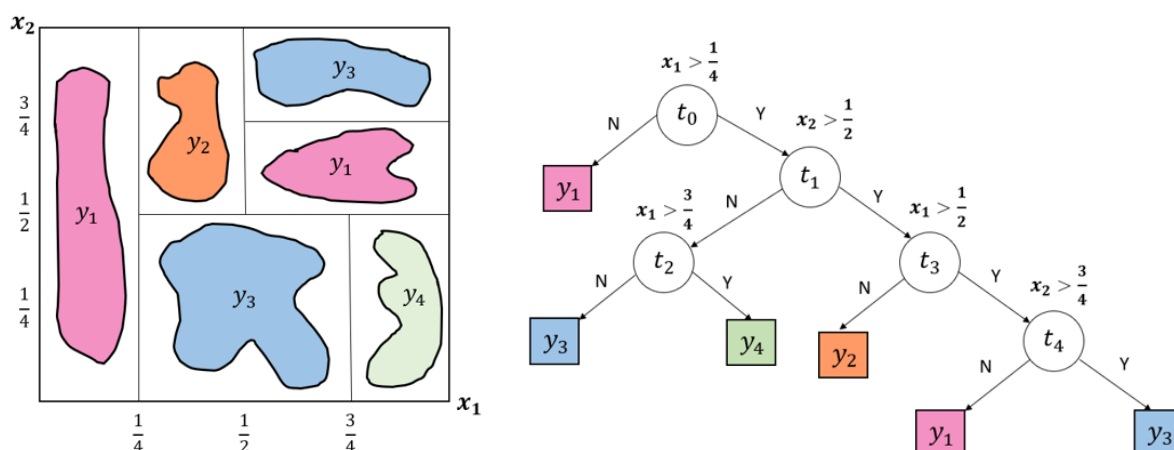


Figure 2: Decision tree classification. (a) OBCT divides the feature space into sub-regions. (b) A set of questions to be asked for classification.

**Splitting criteria.** Every binary split of a node generates two descendant nodes. For the tree growing methodology, from the root node down to the leaves, to make sense,

every split must generate subsets that are more “class homogeneous” compared to the ancestor’s subset. This means that the training feature vectors in each one of the new subsets show a higher preference for specific class(es), whereas data in ancestors are more equally distributed among the classes. The goal, therefore, is to define a measure that quantifies node impurity and split the node so that the overall impurity of the descendant nodes is optimally decreased concerning the ancestor node’s impurity.

We consider two different criteria for evaluating the impurity of a given node  $t$ : *Gini impurity* and *Entropy impurity*.

- *Gini impurity*. It’s given by

$$I(t) = 1 - \sum_{k=1}^K p(y_k|t)^2, \quad (1)$$

where  $p(y_k|t)$  is the probability of a vector in node  $t$  belongs to the class  $y_k$ .  $K$  is the total number of classes in the current node. If the node  $t$  contains only one class,  $I(t) = 0$ . Consider a two-class problem, Gini impurity gets its maximum value when the two classes have the same probability, i.e.,  $I(t) = 1 - (0.5^2 + 0.5^2) = 0.5$ .

- *Entropy impurity*. It’s given by

$$I(t) = - \sum_{k=1}^K p(y_k|t) \cdot \log_2 p(y_k|t), \quad (2)$$

where  $\log_2$  is the logarithm with base 2. This is nothing else than the entropy associated with the data contained in node  $t$ , known from Shannon’s Information Theory. Similar to Gini, Entropy impurity gets 0 if the node  $t$  contains only one class. It gets its maximum value when the probability of the two classes is the same, i.e.,  $I(t)_{max} = -(0.5 \cdot \log_2 0.5 + 0.5 \cdot \log_2 0.5) = 1$ .

Computationally, entropy is more complex since it makes use of logarithms and consequently, the calculation of the Gini index will be faster. In practice, Gini impurity is more widely used (e.g., the default splitting criterion in the scikit-learn library is Gini).

Let us denote the splitted nodes by  $tY$  and  $tN$  according to the “Yes” or “No” answer to the single question adopted for the node  $t$ , also referred as the ancestor node. Let  $N_t$  denotes the data contained in  $t$ . The descendant nodes are associated with two new subsets, that is,  $N_{tY}$ ,  $N_{tN}$ , respectively. The decrease in node impurity is defined as

$$\Delta I(t) = I(t) - \frac{N_{tY}}{N_t} I(tY) - \frac{N_{tN}}{N_t} I(tN),$$

where  $I(tY)$ ,  $I(tN)$  are the impurities of the  $tY$  and  $tN$  nodes, respectively. The goal is to adopt, from the set of candidate questions, the one that performs the split leading to the highest decrease of impurity.

**Stopping rule.** The natural question that now arises is when one decides to stop splitting a node and declares it as a leaf of the tree. A possibility is to adopt a threshold  $T$  and stop splitting if the maximum value of  $\Delta I(t)$  over all possible splits, is less than  $T$ . Other alternatives are to stop splitting either if the subset  $N_t$  is small enough or the data in node  $t$  is pure (i.e., all data samples in node  $t$  belong to the same class).

A critical factor in designing a decision tree is its size. As was the case with the MLPs, the size of a tree must be large enough but not too large; otherwise, it tends to learn the particular details of the training set and exhibits poor generalization performance. Experience has shown that the use of a threshold value for the impurity decreases as the stop-splitting rule does not lead to trees of the right size, because it usually stops tree growing either too early or too late. The most commonly used approach is to grow a tree up to a large size first and then prune nodes according to a pruning criterion.

A drawback associated with tree classifiers is their high variance. In practice, it is common that a small change in the training data set results in a very different tree. The reason for this lies in the hierarchical nature of the tree classifiers. In the next section, random forests are introduced for overcoming such limitations.

### 3 Random forests

A common strategy to overcome the high variance of decision tree classifiers is to combine them. Thus, one can exploit their advantages to reach an overall better performance than could be achieved by using each of them separately. An important observation that justifies such an approach is the following. From the different (candidate) decision trees we design to choose the one that results in the best performance (i.e., the highest classification accuracy). However, different trees may fail (to classify correctly) on different data distributions. That is, even the “best” decision tree can fail on datasets that other trees succeed on. To this end, the random forest algorithm is proposed.

Random forest, like its name implies, consists of a large number of individual decision trees that operate as an ensemble. Each tree in the random forest spits out a class prediction and the class with the most votes becomes the model’s prediction. Figure 3 gives an illustration of the random forest model.

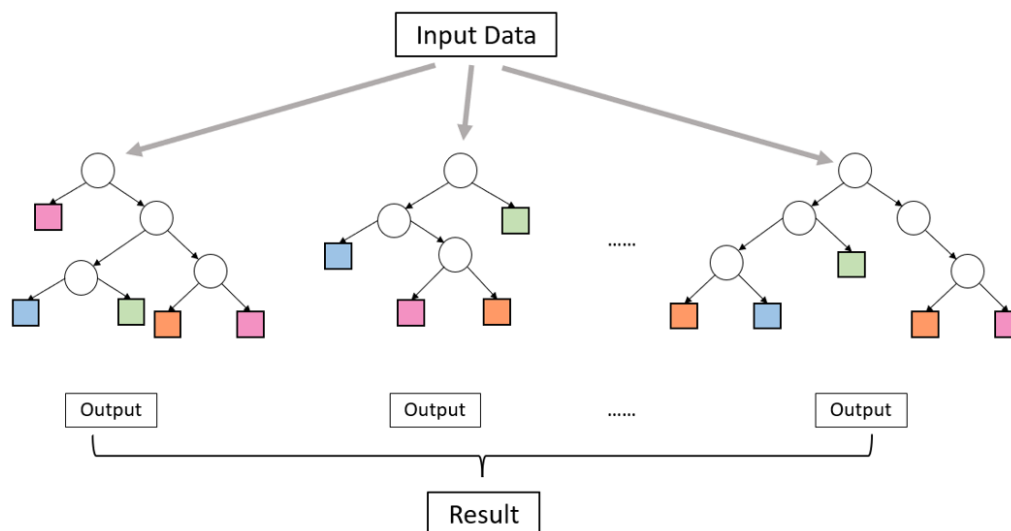


Figure 3: A random forest model.

The fundamental concept behind random forest is a simple but powerful one: the wisdom of crowds. The reason that the random forest model works so well is: *a large number of relatively uncorrelated classifiers (trees) operating as a committee will outperform any of the individual constituent classifiers*. Two aspects are involved to ensure the

independence among the various tree classifiers in the whole forest:

- **Bagging** (also referred to as *bootstrap aggregation*). It is a technique that can reduce variance and improve the generalization error performance. The basic idea is to create a number (say  $B$ ) of variants,  $X_1, X_2, X_3, \dots, X_B$ , of the training set, by uniformly sampling from the original dataset  $X$  with replacement<sup>1</sup>. For each of the training set variants  $X_i$ , a tree  $T_i$  is constructed. The final decision is in favor of the class predicted by the majority of the sub-classifiers,  $T_i (i = 1, 2, 3, \dots, B)$ . Note that with replacement, we are not splitting the training data into smaller chunks and training each tree on a different chunk. Rather, if we have a sample of size  $N$ , we are still feeding each tree a training set of size  $N$  (unless specified otherwise). But instead of the original training data, we take a random sample of size  $N$  with a certain level of data repetitiveness. For instance, if our training data was  $[1, 2, 3, 4, 5, 6]$  then we might give one of our trees the following list  $[1, 2, 2, 3, 6, 6]$ . Notice that both lists are of length six and that “2” and “6” are both repeated in the randomly selected training data we give to our tree (because we sample with replacement).
- **Random feature selection.** In a normal decision tree, when it is time to split a node, we consider every possible feature and pick the one that produces the most separation between the observations in the left node vs. those in the right node using impurity measures. In contrast, each tree in a random forest can pick only from a random subset of features. This forces more variation amongst the trees in the model and ultimately results in lower correlation across trees and more diversification. Figure 4 gives an illustration of random feature selection.

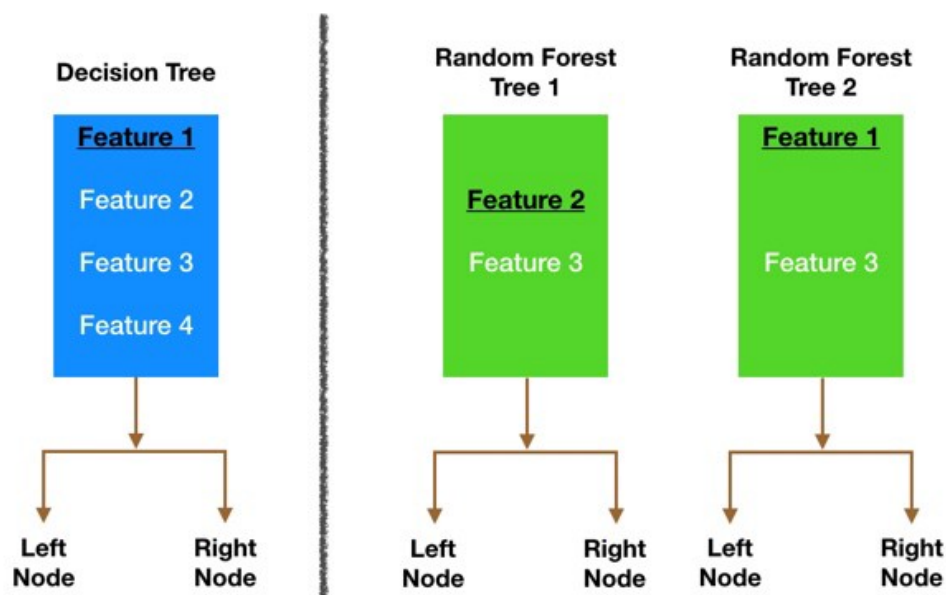


Figure 4: Random feature selection.

Finally, it must be stated that there are close similarities between the decision trees (random forests) and the neural network classifiers. Both aim at forming complex decision boundaries in the feature space. A major difference lies in the way decisions are

<sup>1</sup>It means that we can select the same value multiple times.

made. Decision trees (random forests) employ a hierarchically structured decision function sequentially. In contrast, neural networks utilize a set of soft (not final) decisions in a parallel fashion. Furthermore, their training is performed via different philosophies. However, despite their differences, it has been shown that linear tree classifiers (with a linear splitting criterion) can be adequately mapped to an MLP structure. So far, from the performance point of view, comparative studies seem to give an advantage to the MLPs concerning the classification error, and an advantage to the decision trees (random forests) concerning the required training time.