

Decision tree learning

Decision tree learning is one of the predictive modelling approaches used in statistics, data mining and machine learning. It uses a decision tree (as a predictive model) to go from observations about an item (represented in the branches) to conclusions about the item's target value (represented in the leaves). Tree models where the target variable can take a discrete set of values are called **classification trees**; in these tree structures, leaves represent class labels and branches represent conjunctions of features that lead to those class labels. Decision trees where the target variable can take continuous values (typically real numbers) are called **regression trees**. Decision trees are among the most popular machine learning algorithms given their intelligibility and simplicity.^{[1][2]}

In decision analysis, a decision tree can be used to visually and explicitly represent decisions and decision making. In data mining, a decision tree describes data (but the resulting classification tree can be an input for decision making). This page deals with decision trees in data mining.

Contents

General

Decision tree types

Metrics

Gini impurity

Information gain

Variance reduction

Measure of "goodness"

Uses

Advantages

Limitations

Implementations

Extensions

Decision graphs

Alternative search methods

See also

References

Further reading

External links

General

Decision tree learning is a method commonly used in data mining.^[3] The goal is to create a model that predicts the value of a target variable based on several input variables.

A decision tree is a simple representation for classifying examples. For this section, assume that all of the input features have finite discrete domains, and there is a single target feature called the "classification". Each element of the domain of the classification is called a *class*. A decision tree or a classification tree is a tree in which each internal (non-leaf) node is labeled with an input feature. The arcs coming from a node labeled with an input feature are labeled with each of the possible values of the target feature or the arc leads to a subordinate decision node on a different input feature. Each leaf of the tree is labeled with a class or a probability distribution over the classes, signifying that the data set has been classified by the tree into either a specific class, or into a particular probability distribution (which, if the decision tree is well-constructed, is skewed towards certain subsets of classes).

A tree is built by splitting the source set, constituting the root node of the tree, into subsets—which constitute the successor children. The splitting is based on a set of splitting rules based on classification features.^[4] This process is repeated on each derived subset in a recursive manner called recursive partitioning. The recursion is completed when the subset at a node has all the same values of the target variable, or when splitting no longer adds value to the predictions. This process of *top-down induction of decision trees* (TDIDT)^[5] is an example of a greedy algorithm, and it is by far the most common strategy for learning decision trees from data.

In data mining, decision trees can be described also as the combination of mathematical and computational techniques to aid the description, categorization and generalization of a given set of data.

Data comes in records of the form:

$$(\mathbf{x}, Y) = (x_1, x_2, x_3, \dots, x_k, Y)$$

The dependent variable, Y , is the target variable that we are trying to understand, classify or generalize. The vector \mathbf{x} is composed of the features, x_1, x_2, x_3 etc., that are used for that task.

Decision tree types

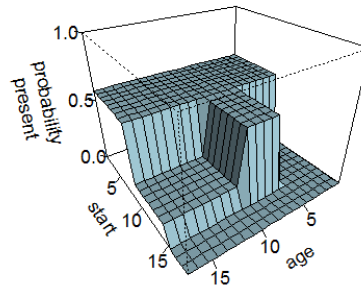
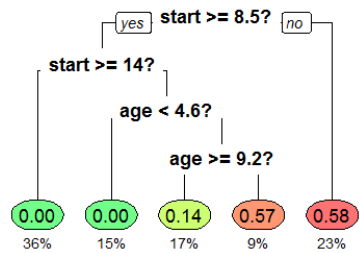
Decision trees used in data mining are of two main types:

- **Classification tree** analysis is when the predicted outcome is the class (discrete) to which the data belongs.
- **Regression tree** analysis is when the predicted outcome can be considered a real number (e.g. the price of a house, or a patient's length of stay in a hospital).

The term **Classification And Regression Tree (CART)** analysis is an umbrella term used to refer to both of the above procedures, first introduced by Breiman et al. in 1984.^[6] Trees used for regression and trees used for classification have some similarities - but also some differences, such as the procedure used to determine where to split.^[6]

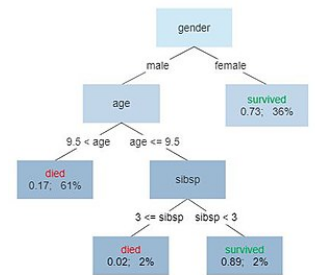
Some techniques, often called *ensemble* methods, construct more than one decision tree:

- **Boosted trees** Incrementally building an ensemble by training each new instance to emphasize the training instances previously mis-modeled. A typical example is AdaBoost. These can be used for

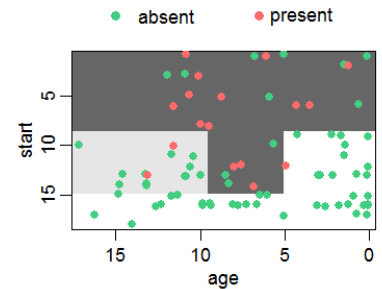


An example tree which estimates the probability of kyphosis after surgery, given the age of the patient and the vertebra at which surgery was started. The same tree is shown in three different ways. **Left** The colored leaves show the probability of kyphosis after surgery, and percentage of patients in the leaf. **Middle** The tree as a perspective plot. **Right** Aerial view of the middle plot. The probability of kyphosis after surgery is higher in the darker areas. (Note: The treatment of kyphosis has advanced considerably since this rather small set of data was collected.)

Survival of passengers on the Titanic



A tree showing survival of passengers on the Titanic ("sibsp" is the number of spouses or siblings aboard). The figures under the leaves show the probability of survival and the percentage of observations in the leaf. Summarizing: Your chances of survival were good if you were (i) a female or (ii) a male younger than 9.5 years with strictly less than 3 siblings.



regression-type and classification-type problems.^{[7][8]}

- **Bootstrap aggregated** (or bagged) decision trees, an early ensemble method, builds multiple decision trees by repeatedly resampling training data with replacement, and voting the trees for a consensus prediction.^[9]
 - A **random forest** classifier is a specific type of bootstrap aggregating
- **Rotation forest** – in which every decision tree is trained by first applying principal component analysis (PCA) on a random subset of the input features.^[10]

A special case of a decision tree is a decision list,^[11] which is a one-sided decision tree, so that every internal node has exactly 1 leaf node and exactly 1 internal node as a child (except for the bottommost node, whose only child is a single leaf node). While less expressive, decision lists are arguably easier to understand than general decision trees due to their added sparsity, permit non-greedy learning methods^[12] and monotonic constraints to be imposed.^[13]

Notable decision tree algorithms include:

- ID3 (Iterative Dichotomiser 3)
- C4.5 (successor of ID3)
- CART (Classification And Regression Tree)^[6]
- Chi-square automatic interaction detection (CHAID). Performs multi-level splits when computing classification trees.^[14]
- MARS: extends decision trees to handle numerical data better.

- Conditional Inference Trees. Statistics-based approach that uses non-parametric tests as splitting criteria, corrected for multiple testing to avoid overfitting. This approach results in unbiased predictor selection and does not require pruning.^{[15][16]}

ID3 and CART were invented independently at around the same time (between 1970 and 1980), yet follow a similar approach for learning a decision tree from training tuples.

It has also been proposed to leverage concepts of fuzzy set theory for the definition of a special version of decision tree, known as Fuzzy Decision Tree (FDT).^[17] In this type of fuzzy classification, generally an input vector \mathbf{x} is associated with multiple classes, each with a different confidence value. Boosted ensembles of FDTs have been recently investigated as well, and they have shown performances comparable to those of other very efficient fuzzy classifiers.^[18]

Metrics

Algorithms for constructing decision trees usually work top-down, by choosing a variable at each step that best splits the set of items.^[19] Different algorithms use different metrics for measuring "best". These generally measure the homogeneity of the target variable within the subsets. Some examples are given below. These metrics are applied to each candidate subset, and the resulting values are combined (e.g., averaged) to provide a measure of the quality of the split.

Gini impurity

Used by the CART (classification and regression tree) algorithm for classification trees, Gini impurity is a measure of how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset. The Gini impurity can be computed by summing the probability p_i of an item with label i being chosen times the probability $\sum_{k \neq i} p_k = 1 - p_i$ of a mistake in categorizing that item. It reaches its minimum (zero) when all cases in the node fall into a single target category.

The Gini impurity is also an information theoretic measure and corresponds to Tsallis Entropy with deformation coefficient $q = 2$, which in physics is associated with the lack of information in out-of-equilibrium, non-extensive, dissipative and quantum systems. For the limit $q \rightarrow 1$ one recovers the usual Boltzmann-Gibbs or Shannon entropy. In this sense, the Gini impurity is but a variation of the usual entropy measure for decision trees.

To compute Gini impurity for a set of items with J classes, suppose $i \in \{1, 2, \dots, J\}$, and let p_i be the fraction of items labeled with class i in the set.

$$I_G(p) = \sum_{i=1}^J \left(p_i \sum_{k \neq i} p_k \right) = \sum_{i=1}^J p_i (1 - p_i) = \sum_{i=1}^J (p_i - p_i^2) = \sum_{i=1}^J p_i - \sum_{i=1}^J p_i^2 = 1 - \sum_{i=1}^J p_i^2$$

Information gain

Used by the ID3, C4.5 and C5.0 tree-generation algorithms. Information gain is based on the concept of entropy and information content from information theory.

Entropy is defined as below

$$H(T) = I_E(p_1, p_2, \dots, p_J) = - \sum_{i=1}^J p_i \log_2 p_i$$

where p_1, p_2, \dots are fractions that add up to 1 and represent the percentage of each class present in the child node that results from a split in the tree.^[20]

$$\begin{aligned} \text{Information Gain} &= \text{Entropy (parent)} - \text{Sum of Entropy (Children)} \\ \overbrace{IG(T, a)} &= \overbrace{H(T)} - \overbrace{H(T|a)} \\ &= - \sum_{i=1}^J p_i \log_2 p_i - \sum_{i=1}^J - \text{Pr}(i|a) \log_2 \text{Pr}(i|a) \end{aligned}$$

Averaging over the possible values of A ,

$$\begin{aligned} \underbrace{\text{Expected Information Gain}}_{E_A(IG(T, a))} &= \underbrace{\text{Mutual Information between T and A}}_{I(T; A)} = \underbrace{\text{Entropy (parent)}}_{H(T)} - \underbrace{\text{Weighted Sum of Entropy (Children)}}_{H(T|A)} \\ &= - \sum_{i=1}^J p_i \log_2 p_i - \sum_a p(a) \sum_{i=1}^J -\Pr(i|a) \log_2 \Pr(i|a) \end{aligned}$$

That is, the expected information gain is the mutual information, meaning that on average, the reduction in the entropy of T is the mutual information.

Information gain is used to decide which feature to split on at each step in building the tree. Simplicity is best, so we want to keep our tree small. To do so, at each step we should choose the split that results in the purest daughter nodes. A commonly used measure of purity is called information which is measured in bits. For each node of the tree, the information value "represents the expected amount of information that would be needed to specify whether a new instance should be classified yes or no, given that the example reached that node".^[20]

Consider an example data set with four attributes: *outlook* (sunny, overcast, rainy), *temperature* (hot, mild, cool), *humidity* (high, normal), and *windy* (true, false), with a binary (yes or no) target variable, *play*, and 14 data points. To construct a decision tree on this data, we need to compare the information gain of each of four trees, each split on one of the four features. The split with the highest information gain will be taken as the first split and the process will continue until all children nodes are pure, or until the information gain is 0.

To find the information gain of the split using *windy*, we must first calculate the information in the data before the split. The original data contained nine yes's and five no's.

$$I_E([9, 5]) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

The split using the feature *windy* results in two children nodes, one for a *windy* value of true and one for a *windy* value of false. In this data set, there are six data points with a true *windy* value, three of which have a *play* (where *play* is the target variable) value of yes and three with a *play* value of no. The eight remaining data points with a *windy* value of false contain two no's and six yes's. The information of the *windy*=true node is calculated using the entropy equation above. Since there is an equal number of yes's and no's in this node, we have

$$I_E([3, 3]) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

For the node where *windy*=false there were eight data points, six yes's and two no's. Thus we have

$$I_E([6, 2]) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.81$$

To find the information of the split, we take the weighted average of these two numbers based on how many observations fell into which node.

$$I_E([3, 3], [6, 2]) = I_E(\text{windy or not}) = \frac{6}{14} \cdot 1 + \frac{8}{14} \cdot 0.81 = 0.89$$

Now we can calculate the information gain achieved by splitting on the *windy* feature.

$$IG(\text{windy}) = I_E([9, 5]) - I_E([3, 3], [6, 2]) = 0.94 - 0.89 = 0.05$$

To build the tree, the information gain of each possible first split would need to be calculated. The best first split is the one that provides the most information gain. This process is repeated for each impure node until the tree is complete. This example is adapted from the example appearing in Witten et al.^[20]

Variance reduction

Introduced in CART,^[6] variance reduction is often employed in cases where the target variable is continuous (regression tree), meaning that use of many other metrics would first require discretization before being applied. The variance reduction of a node N is defined as the total reduction of the variance of the target variable Y due to the split at this node:

$$I_V(N) = \frac{1}{|S|^2} \sum_{i \in S} \sum_{j \in S} \frac{1}{2} (y_i - y_j)^2 - \left(\frac{1}{|S_t|^2} \sum_{i \in S_t} \sum_{j \in S_t} \frac{1}{2} (y_i - y_j)^2 + \frac{1}{|S_f|^2} \sum_{i \in S_f} \sum_{j \in S_f} \frac{1}{2} (y_i - y_j)^2 \right)$$

where S , S_t , and S_f are the set of presplit sample indices, set of sample indices for which the split test is true, and set of sample indices for which the split test is false, respectively. Each of the above summands are indeed variance estimates, though, written in a form without directly referring to the mean.

Measure of "goodness"

Used by CART in 1984,^[21] the measure of "goodness" is a function that seeks to optimize the balance of a candidate split's capacity to create pure children with its capacity to create equally-sized children. This process is repeated for each impure node until the tree is complete. The function $\phi(s|t)$, where s is a candidate split at node t , is defined as below

$$\phi(s|t) = 2P_L P_R \sum_{j=1}^{ClassCount} |P(j|t_L) - P(j|t_R)|$$

where t_L and t_R are the left and right children of node t using split s , respectively; P_L and P_R are the proportions of records in t in t_L and t_R , respectively; and $P(j|t_L)$ and $P(j|t_R)$ are the proportions of class j records in t_L and t_R , respectively.

Consider an example data set with three attributes: *savings*(low, medium, high), *assets*(low, medium, high), *income*(numerical value), and a binary target variable *credit risk*(good, bad) and 8 data points.^[21] The full data is presented in the table below. To start a decision tree, we will calculate the maximum value of $\phi(s|t)$ using each feature to find which one will split the root node. This process will continue until all children are pure or all $\phi(s|t)$ values are below a set threshold.

Customer	Savings	Assets	Income (\$1000s)	Credit Risk
1	Medium	High	75	Good
2	Low	Low	50	Bad
3	High	Medium	25	Bad
4	Medium	Medium	50	Good
5	Low	Medium	100	Good
6	High	High	25	Good
7	Low	Low	25	Bad
8	Medium	Medium	75	Good

To find $\phi(s|t)$ of the feature *savings*, we need to note the quantity of each value. The original data contained three low's, three medium's, and two high's. Out of the low's, one had a good *credit risk* while out of the medium's and high's, 4 had a good *credit risk*. Assume a candidate split s such that records with a low *savings* will be put in the left child and all other records will be put into the right child.

$$\phi(s|root) = 2 \frac{3}{8} * \frac{5}{8} * (|\frac{1}{3} - \frac{4}{5}| + |\frac{2}{3} - \frac{1}{5}|) = 0.44$$

To build the tree, the "goodness" of all candidate splits for the root node need to be calculated. The candidate with the maximum value will split the root node, and the process will continue for each impure node until the tree is complete.

Compared to other metrics such as information gain, the measure of "goodness" will attempt to create a more balanced tree, leading to more-consistent decision time. However, it sacrifices some priority for creating pure children which can lead to additional splits that are not present with other metrics.

Uses

Advantages

Amongst other data mining methods, decision trees have various advantages:

- **Simple to understand and interpret.** People are able to understand decision tree models after a brief explanation. Trees can also be displayed graphically in a way that is easy for non-experts to interpret.^[22]
- **Able to handle both numerical and categorical data.**^[22] Other techniques are usually specialized in analyzing datasets that have only one type of variable. (For example, relation rules can be used only with nominal variables while neural networks can be used only with numerical variables or categoricals converted to 0-1 values.) Early decision trees were only capable of handling categorical variables, but more recent versions, such as C4.5, do not have this limitation.^[2]
- **Requires little data preparation.** Other techniques often require data normalization. Since trees can handle qualitative predictors, there is no need to create dummy variables.^[22]
- **Uses a white box or open-box**^[2] **model.** If a given situation is observable in a model the explanation for the condition is easily explained by boolean logic. By contrast, in a black box model, the explanation for the results is typically difficult to understand, for example with an artificial neural network.
- **Possible to validate a model using statistical tests.** That makes it possible to account for the reliability of the model.
- Non-statistical approach that makes no assumptions of the training data or prediction residuals; e.g., no distributional, independence, or constant variance assumptions
- **Performs well with large datasets.** Large amounts of data can be analyzed using standard computing resources in reasonable time.
- **Mirrors human decision making more closely than other approaches.**^[22] This could be useful when modeling human decisions/behavior.
- **Robust against co-linearity, particularly boosting**
- **In built feature selection.** Additional irrelevant feature will be less used so that they can be removed on subsequent runs. The hierarchy of attributes in a decision tree reflects the importance of attributes.^[23] It means that the features on top are the most informative.^[24]
- **Decision trees can approximate any Boolean function eq. XOR.**^[25]

Limitations

- Trees can be very non-robust. A small change in the training data can result in a large change in the tree and consequently the final predictions.^[22]
- The problem of learning an optimal decision tree is known to be NP-complete under several aspects of optimality and even for simple concepts.^{[26][27]} Consequently, practical decision-tree learning algorithms are based on heuristics such as the greedy algorithm where locally optimal decisions are made at each node. Such algorithms cannot guarantee to return the globally optimal decision tree. To reduce the greedy effect of local optimality, some methods such as the dual information distance (DID) tree were proposed.^[28]
- Decision-tree learners can create over-complex trees that do not generalize well from the training data. (This is known as overfitting.^[29]) Mechanisms such as pruning are necessary to avoid this problem (with the exception of some algorithms such as the Conditional Inference approach, that does not require pruning).^{[15][16]}
- The average depth of the tree that is defined by the number of nodes or tests till classification is not guaranteed to be minimal or small under various splitting criteria.^[30]
- For data including categorical variables with different numbers of levels, information gain in decision trees is biased in favor of attributes with more levels.^[31] However, the issue of biased predictor selection is avoided by the Conditional Inference approach,^[15] a two-stage approach,^[32] or adaptive leave-one-out feature selection.^[33]

Implementations

Many data mining software packages provide implementations of one or more decision tree algorithms.

Examples include

- Salford Systems CART (which licensed the proprietary code of the original CART authors),^[6]
- IBM SPSS Modeler,
- RapidMiner,
- SAS Enterprise Miner,
- Matlab,

- [R](#) (an open-source software environment for statistical computing, which includes several CART implementations such as [rpart](#), [party](#) and [randomForest](#) packages),
- [Weka](#) (a free and open-source data-mining suite, contains many decision tree algorithms),
- [Orange](#),
- [KNIME](#),
- [Microsoft SQL Server](#) [1] (<https://technet.microsoft.com/en-us/library/cc645868.aspx>), and
- [scikit-learn](#) (a free and open-source machine learning library for the [Python](#) programming language).

Extensions

Decision graphs

In a decision tree, all paths from the root node to the leaf node proceed by way of conjunction, or *AND*. In a decision graph, it is possible to use disjunctions (ORs) to join two more paths together using [minimum message length](#) (MML).^[34] Decision graphs have been further extended to allow for previously unstated new attributes to be learnt dynamically and used at different places within the graph.^[35] The more general coding scheme results in better predictive accuracy and log-loss probabilistic scoring. In general, decision graphs infer models with fewer leaves than decision trees.

Alternative search methods

Evolutionary algorithms have been used to avoid local optimal decisions and search the decision tree space with little *a priori* bias.^{[36][37]}

It is also possible for a tree to be sampled using [MCMC](#).^[38]

The tree can be searched for in a bottom-up fashion.^[39] Or several trees can be constructed parallelly to reduce the expected number of tests till classification.^[30]

See also

- [Decision tree pruning](#)
- [Binary decision diagram](#)
- [CHAID](#)
- [CART](#)
- [ID3 algorithm](#)
- [C4.5 algorithm](#)
- [Decision stumps](#), used in e.g. [AdaBoosting](#)
- [Decision list](#)
- [Incremental decision tree](#)
- [Alternating decision tree](#)
- [Structured data analysis \(statistics\)](#)
- [Logistic model tree](#)
- [Hierarchical clustering](#)

References

- Wu, Xindong; Kumar, Vipin; Ross Quinlan, J.; Ghosh, Joydeep; Yang, Qiang; Motoda, Hiroshi; McLachlan, Geoffrey J.; Ng, Angus; Liu, Bing; Yu, Philip S.; Zhou, Zhi-Hua (2008-01-01). "Top 10 algorithms in data mining". *Knowledge and Information Systems*. **14** (1): 1–37. doi:10.1007/s10115-007-0114-2 (<https://doi.org/10.1007/s10115-007-0114-2>). ISSN 0219-3116 (<http://www.worldcat.org/issn/0219-3116>). S2CID 2367747 (<https://api.semanticscholar.org/CorpusID:2367747>).
- Piryonesi S. Madeh; El-Diraby Tamer E. (2020-03-01). "Data Analytics in Asset Management: Cost-Effective Prediction of the Pavement Condition Index". *Journal of Infrastructure Systems*. **26** (1): 04019036. doi:10.1061/(ASCE)IS.1943-555X.0000512 ([https://doi.org/10.1061/\(ASCE\)IS.1943-555X.0000512](https://doi.org/10.1061/(ASCE)IS.1943-555X.0000512)).
- Rokach, Lior; Maimon, O. (2008). *Data mining with decision trees: theory and applications*. World Scientific Pub Co Inc. ISBN 978-9812771711.
- Shalev-Shwartz, Shai; Ben-David, Shai (2014). "18. Decision Trees". *Understanding Machine Learning* (<http://www.cs.huji.ac.il/~shais/UnderstandingMachineLearning>). Cambridge University Press.
- Quinlan, J. R. (1986). "Induction of decision trees" (<https://link.springer.com/content/pdf/10.1007/BF00116251.pdf>) (PDF). *Machine Learning*. **1**: 81–106. doi:10.1007/BF00116251 (<https://doi.org/10.1007/BF00116251>). S2CID 189902138 (<https://api.semanticscholar.org/CorpusID:189902138>).
- Breiman, Leo; Friedman, J. H.; Olshen, R. A.; Stone, C. J. (1984). *Classification and regression trees*. Monterey, CA: Wadsworth & Brooks/Cole Advanced Books & Software. ISBN 978-0-412-04841-8.
- Friedman, J. H. (1999). *Stochastic gradient boosting* ([https://astro.temple.edu/~msobel/courses_files/StochasticBoosting\(gradient\).pdf](https://astro.temple.edu/~msobel/courses_files/StochasticBoosting(gradient).pdf)). Stanford University.
- Hastie, T., Tibshirani, R., Friedman, J. H. (2001). *The elements of statistical learning : Data mining, inference, and prediction*. New York: Springer Verlag.

9. Breiman, L. (1996). "Bagging Predictors" (<https://doi.org/10.1007%2FBF00058655>). *Machine Learning*. **24** (2): 123–140. doi:10.1007/BF00058655 (<https://doi.org/10.1007%2FBF00058655>).
10. Rodriguez, J. J.; Kuncheva, L. I.; Alonso, C. J. (2006). "Rotation forest: A new classifier ensemble method". *IEEE Transactions on Pattern Analysis and Machine Intelligence*. **28** (10): 1619–1630. CiteSeerX 10.1.1.156.8277 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.156.8277>). doi:10.1109/TPAMI.2006.211 (<https://doi.org/10.1109%2FTPAMI.2006.211>). PMID 16986543 (<https://pubmed.ncbi.nlm.nih.gov/16986543>). S2CID 6847493 (<https://api.semanticscholar.org/CorpusID:6847493>).
11. Rivest, Ron (Nov 1987). "Learning Decision Lists" (<http://people.csail.mit.edu/rivest/pubs/Riv87b.pdf>) (PDF). *Machine Learning*. **3** (2): 229–246. doi:10.1023/A:1022607331053 (<https://doi.org/10.1023%2FA%3A1022607331053>). S2CID 30625841 (<https://api.semanticscholar.org/CorpusID:30625841>).
12. Letham, Ben; Rudin, Cynthia; McCormick, Tyler; Madigan, David (2015). "Interpretable Classifiers Using Rules And Bayesian Analysis: Building A Better Stroke Prediction Model". *Annals of Applied Statistics*. **9** (3): 1350–1371. arXiv:1511.01644 (<https://arxiv.org/abs/1511.01644>). doi:10.1214/15-AOAS848 (<https://doi.org/10.1214%2F15-AOAS848>). S2CID 17699665 (<https://api.semanticscholar.org/CorpusID:17699665>).
13. Wang, Fulton; Rudin, Cynthia (2015). "Falling Rule Lists" (<http://www.jmlr.org/proceedings/papers/v38/wang15a.pdf>) (PDF). *Journal of Machine Learning Research*. **38**.
14. Kass, G. V. (1980). "An exploratory technique for investigating large quantities of categorical data". *Applied Statistics*. **29** (2): 119–127. doi:10.2307/2986296 (<https://doi.org/10.2307%2F2986296>). JSTOR 2986296 (<https://www.jstor.org/stable/2986296>).
15. Hothorn, T.; Hornik, K.; Zeileis, A. (2006). "Unbiased Recursive Partitioning: A Conditional Inference Framework". *Journal of Computational and Graphical Statistics*. **15** (3): 651–674. CiteSeerX 10.1.1.527.2935 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.527.2935>). doi:10.1198/106186006X133933 (<https://doi.org/10.1198%2F106186006X133933>). JSTOR 27594202 (<https://www.jstor.org/stable/27594202>). S2CID 6074128 (<https://api.semanticscholar.org/CorpusID:6074128>).
16. Strobl, C.; Malley, J.; Tutz, G. (2009). "An Introduction to Recursive Partitioning: Rationale, Application and Characteristics of Classification and Regression Trees, Bagging and Random Forests" (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2927982>). *Psychological Methods*. **14** (4): 323–348. doi:10.1037/a0016973 (<https://doi.org/10.1037%2Fa0016973>). PMC 2927982 (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2927982>). PMID 19968396 (<https://pubmed.ncbi.nlm.nih.gov/19968396>).
17. Janikow, C. Z. (1998). "Fuzzy decision trees: issues and methods". *IEEE Transactions on Systems, Man and Cybernetics, Part B (Cybernetics)*. **28** (1): 1–14. doi:10.1109/3477.658573 (<https://doi.org/10.1109%2F3477.658573>). PMID 18255917 (<https://pubmed.ncbi.nlm.nih.gov/18255917>).
18. Barsacchi, M.; Bechini, A.; Marcelloni, F. (2020). "An analysis of boosted ensembles of binary fuzzy decision trees" (<http://www.sciencedirect.com/science/article/pii/S0957417420302608>). *Expert Systems with Applications*. **154**: 113436. doi:10.1016/j.eswa.2020.113436 (<https://doi.org/10.1016%2Fj.eswa.2020.113436>).
19. Rokach, L.; Maimon, O. (2005). "Top-down induction of decision trees classifiers-a survey". *IEEE Transactions on Systems, Man, and Cybernetics - Part C: Applications and Reviews*. **35** (4): 476–487. CiteSeerX 10.1.1.458.7031 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.458.7031>). doi:10.1109/TSMCC.2004.843247 (<https://doi.org/10.1109%2FTSMCC.2004.843247>). S2CID 14808716 (<https://api.semanticscholar.org/CorpusID:14808716>).
20. Witten, Ian; Frank, Eibe; Hall, Mark (2011). *Data Mining* (https://archive.org/details/dataminingpracti00witt_966). Burlington, MA: Morgan Kaufmann. pp. 102 (https://archive.org/details/dataminingpracti00witt_966/page/n136)–103. ISBN 978-0-12-374856-0.
21. Larose, Daniel T.; Larose, Chantal D. (2014). *Discovering knowledge in data: an introduction to data mining*. Hoboken, NJ: John Wiley & Sons, Inc. ISBN 9781118874059.
22. Gareth, James; Witten, Daniela; Hastie, Trevor; Tibshirani, Robert (2015). *An Introduction to Statistical Learning* (<https://archive.org/details/introductiontost00jame>). New York: Springer. pp. 315 (<https://archive.org/details/introductiontost00jame/page/n323>). ISBN 978-1-4614-7137-0.
23. Provost, Foster, 1964- (2013). *Data science for business : [what you need to know about data mining and data-analytic thinking]*. Fawcett, Tom. (1st ed.). Sebastopol, Calif.: O'Reilly. ISBN 978-1-4493-6132-7. OCLC 844460899 (<https://www.worldcat.org/oclc/844460899>).
24. Pirayonesi S. Madeh; El-Diraby Tamer E. (2020-06-01). "Role of Data Analytics in Infrastructure Asset Management: Overcoming Data Size and Quality Problems". *Journal of Transportation Engineering, Part B: Pavements*. **146** (2): 04020022. doi:10.1061/JPEODX.0000175 (<https://doi.org/10.1061%2FJPEODX.0000175>).
25. Mehtaa, Dinesh; Raghavan, Vijay (2002). "Decision tree approximations of Boolean functions". *Theoretical Computer Science*. **270** (1–2): 609–623. doi:10.1016/S0304-3975(01)00011-1 ([https://doi.org/10.1016%2FS0304-3975\(01\)00011-1](https://doi.org/10.1016%2FS0304-3975(01)00011-1)).
26. Hyafil, Laurent; Rivest, RL (1976). "Constructing Optimal Binary Decision Trees is NP-complete". *Information Processing Letters*. **5** (1): 15–17. doi:10.1016/0020-0190(76)90095-8 ([https://doi.org/10.1016%2F0020-0190\(76\)90095-8](https://doi.org/10.1016%2F0020-0190(76)90095-8)).
27. Murthy S. (1998). "Automatic construction of decision trees from data: A multidisciplinary survey" (https://cs.nyu.edu/~roweis/csc2515-2006/readings/murthy_dt.pdf). *Data Mining and Knowledge Discovery*
28. Ben-Gal I. Dana A., Shkolnik N. and Singer (2014). "Efficient Construction of Decision Trees by the Dual Information Distance Method" (<http://www.eng.tau.ac.il/~bengal/DID.pdf>) (PDF). *Quality Technology & Quantitative Management*. **11** (1): 133–147. doi:10.1080/16843703.2014.11673330 (<https://doi.org/10.1080%2F16843703.2014.11673330>). S2CID 7025979 (<https://api.semanticscholar.org/CorpusID:7025979>).
29. *Principles of Data Mining*. 2007. doi:10.1007/978-1-84628-766-4 (<https://doi.org/10.1007%2F978-1-84628-766-4>). ISBN 978-1-84628-766-7.
30. Ben-Gal I. and Trister C. (2015). "Parallel Construction of Decision Trees with Consistently Non Increasing Expected Number of Tests " (<http://www.eng.tau.ac.il/~bengal/Trist.pdf>) (PDF). *Applied Stochastic Models in Business and Industry*, Vol. 31(1) 64–78.

31. Deng, H.; Runger, G.; Tuv, E. (2011). *Bias of importance measures for multi-valued attributes and solutions* (<https://www.researchgate.net/publication/221079908>). Proceedings of the 21st International Conference on Artificial Neural Networks (ICANN). pp. 293–300.
32. Brandmaier, Andreas M.; Oertzen, Timo von; McArdle, John J.; Lindenberger, Ulman (2012). "Structural equation model trees" (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4386908>). *Psychological Methods*. **18** (1): 71–86. doi:10.1037/a0030001 (<https://doi.org/10.1037/a0030001>). hdl:11858/00-001M-0000-0024-EA33-9 (<https://hdl.handle.net/11858/00-001M-0000-0024-EA33-9>). PMC 4386908 (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4386908>). PMID 22984789 (<https://pubmed.ncbi.nlm.nih.gov/22984789>).
33. Painsky, Amichai; Rosset, Saharon (2017). "Cross-Validated Variable Selection in Tree-Based Methods Improves Predictive Performance". *IEEE Transactions on Pattern Analysis and Machine Intelligence*. **39** (11): 2142–2153. arXiv:1512.03444 (<https://arxiv.org/abs/1512.03444>). doi:10.1109/TPAMI.2016.2636831 (<https://doi.org/10.1109/TPAMI.2016.2636831>). PMID 28114007 (<https://pubmed.ncbi.nlm.nih.gov/28114007>). S2CID 5381516 (<https://api.semanticscholar.org/CorpusID:5381516>).
34. "CiteSeerX" (<http://citeseer.ist.psu.edu/oliver93decision.html>).
35. Tan & Dowe (2003) (http://www.csse.monash.edu.au/~dld/Publications/2003/Tan+Dowe2003_MMLDecisionGraphs.pdf)
36. Papagelis, A.; Kalles, D. (2001). "Breeding Decision Trees Using Evolutionary Techniques" (<http://www.gatree.com/wordpress/wp-content/uploads/2010/04/BreedingDecisionTreeUsingEvolutionaryTechniques.pdf>) (PDF). *Proceedings of the Eighteenth International Conference on Machine Learning, June 28–July 1, 2001*. pp. 393–400.
37. Barros, Rodrigo C.; Basgalupp, M. P.; Carvalho, A. C. P. L. F.; Freitas, Alex A. (2012). "A Survey of Evolutionary Algorithms for Decision-Tree Induction". *IEEE Transactions on Systems, Man and Cybernetics*. Part C: Applications and Reviews. **42** (3): 291–312. CiteSeerX 10.1.1.308.9068 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.308.9068>). doi:10.1109/TSMCC.2011.2157494 (<https://doi.org/10.1109/TSMCC.2011.2157494>). S2CID 365692 (<https://api.semanticscholar.org/CorpusID:365692>).
38. Chipman, Hugh A.; George, Edward I.; McCulloch, Robert E. (1998). "Bayesian CART model search". *Journal of the American Statistical Association*. **93** (443): 935–948. CiteSeerX 10.1.1.211.5573 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.211.5573>). doi:10.1080/01621459.1998.10473750 (<https://doi.org/10.1080/01621459.1998.10473750>).
39. Barros, R. C.; Cerri, R.; Jaskowiak, P. A.; Carvalho, A. C. P. L. F. (2011). "A bottom-up oblique decision tree induction algorithm". *Proceedings of the 11th International Conference on Intelligent Systems Design and Applications (ISDA 2011)*. pp. 450–456. doi:10.1109/ISDA.2011.6121697 (<https://doi.org/10.1109/ISDA.2011.6121697>). ISBN 978-1-4577-1676-8. S2CID 15574923 (<https://api.semanticscholar.org/CorpusID:15574923>).

Further reading

- James, Gareth; Witten, Daniela; Hastie, Trevor; Tibshirani, Robert (2017). "Tree-Based Methods" (<https://www-bcf.usc.edu/~gareth/ISL/ISLR%20Seventh%20Printing.pdf#page=317>) (PDF). *An Introduction to Statistical Learning: with Applications in R*. New York: Springer. pp. 303–336. ISBN 978-1-4614-7137-0.

External links

- Evolutionary Learning of Decision Trees in C++ (<https://www.cs.kent.ac.uk/people/staff/mg483/code/evoldectrees/>)
- A very detailed explanation of information gain as splitting criterion (<http://christianherta.de/lehre/dataScience/machineLearning/decision-trees.html>)

Retrieved from "https://en.wikipedia.org/w/index.php?title=Decision_tree_learning&oldid=1003832169"

This page was last edited on 30 January 2021, at 22:24 (UTC).

Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.