# Classification tree

* A recursive two-way partition (or branches) of predictors.

nodes

Class identifier.

B
0·65
100%

Yes — Concave_points_mean ≥ 0.56 — No

Splitting rules
(branches)

B
0.92
67%
← Proportion of the reference class.

radius_mean ≥ 15

M
0.09
33%
Percentage of observations in node.

M
0.40
5%

B
0.96
62%

terminal nodes (or leaves)

⇒ Colors indicate class, with a darker color indicating lower impurity.

Pure ←————— Impure —————→ Pure

100% of one label.

75% of one label

50% of each label.

75% of one label

100% of one label

## Notations

* $s$ = the value to split in variable $j$.

* $A_L = \{\{y_i, x_i\} : x_{ij} < s\}$ and

  $A_R = \{\{y_i, x_i\} : x_{ij} \geqslant s\}$.

* $P_{kA}$ = the proportion of observations in class $k$ in a set $A$.

## Impurity metrics

Why we need? $\Rightarrow$ An algorithm to find the optimal $\{j^*, s^*\}$
(split and variable).

selected predictor

selected threshold.

Gini index: $f_{Gini}(A) = \sum\limits_{k=1}^{k} P_{kA}(1 - P_{kA}) = 1 - \sum\limits_{k=1}^{k} P_{kA}^2$

entropy index: $f_{entropy}(A) = -\sum\limits_{k=1}^{k} P_{kA} \log(P_{kA})$

Controls how a decision tree decides where to split the data.

* smaller values of the impurity index means higher purity.

Overall impurity $= \dfrac{|A_L|}{|A_L| + |A_R|} f(A_L) + \dfrac{|A_R|}{|A_L| + |A_R|} f(A_R)$

## Objective.

$$[j^*, s^*] = \underset{j \in \{1,\ldots, p\},\, s \in \mathbb{R}}{\arg\min} \text{Overall impurity } (A_L, A_R).$$

$\Rightarrow$ We have to repeat the process until, we reach the (stopping rule.)

* **minsplit:** the minimum number of observations in any non-terminal node.

* **minbucket:** the minimum number of observations allowed in a terminal node.

* **cp:** complexity parameter — minimum difference between impurity values required to continue splitting.

rel error = in-sample error (always decreases with more split).
Xerror = the cross-validation error.
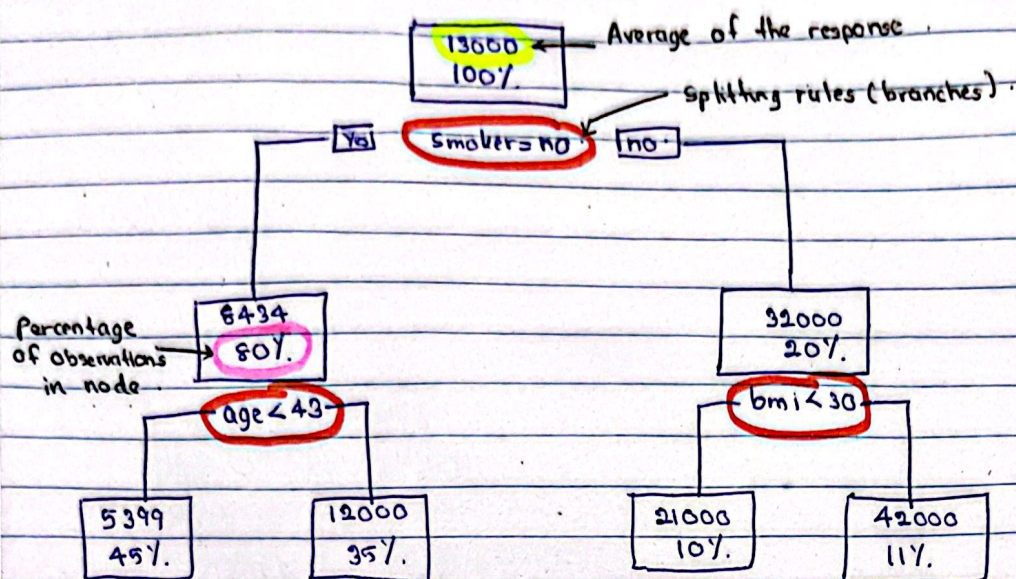Xstd = the standard deviation of the cross-validation error.

## Accuracy measures

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN}.$$

$$\text{Balanced accuracy} = \frac{TPR + FPR}{2} = \frac{TP}{2(TP+FN)} + \frac{FP}{2(FP+TN)}$$

$$\left.\begin{array}{l} \text{Cohen's kappa} \\ \text{coefficient} \end{array}\right\} = \frac{2 \times (TP \times TN - FN \times FP)}{(TP+FP)(TN+FP) + (TP+FN)(TN+FN)}.$$

# Regression tree



Diagram labels:
- 13000 / 100% — Average of the response
- Yes — Smoker = no — no — Splitting rules (branches)
- 6434 / 80% — Percentage of observations in node
- 92000 / 20%
- age < 43
- bmi < 30
- 5399 / 45%
- 12000 / 35%
- 21000 / 10%
- 42000 / 11%

* Colors indicate class, with a darker color indicating lower impurity.

* The prediction is the average of the response in the terminal node.

* For a binary variable, $A_L$ and $A_R$ is based on its class.
* Then find the average responses, $\bar{y}_L$ and $\bar{y}_R$, for $A_L$ and $A_R$, respectively.

* We can use the residual sum of squares to measure the impurity.
$$RSS(A_L) = \sum_{i \in A_L} (y_i - \bar{y}_L)^2.$$

* The overall impurity is $RSS(A_L) + RSS(A_R)$.

## Bagging (More on next tutorial)
* Multiple subsets are created from the original data set with equal tuples, selecting observations with replacement.

* A base model is created on each of these subsets.

* Each model is learned in parallel with each training set and independent of each other.
* The final predictions are determined by combining the predictions from all the models.