- Decision Trees

→ Decision trees can be applied to both regression & classification problems.

① Regression Trees:-



Step 1:- Divide predictor space - i.e set of possible values for $X_1, X_2, \dots X_p$ - into $J$ distinct and non-overlapping regions, $R_1, R_2 \dots R_J$

Step 2:- for every observation that falls in to $R_j$, we make same prediction ( mean of responses in that region)

↳ How to construct regions $R_1, R_2, \dots R_J$?

○ Goal:- find $R_1, \dots R_J$ regions that minimize the RSS given by

$$\sum_{J=1}^{J} \sum_{i \in R_j} \left( y_i - \hat{y}_{R_j} \right)^2$$

○ Top-down, Greedy approach :(Recursive binary Splitting) ↳cause at each step of the tree-building process, the best split is made at the particular step, rather than looking ahead & picking a split that will lead to a better tree in future step

⇒ Stopping Criterion:

process of minimizing RSS continues till a

no region contains

stopping Criterion ( Ex: ≥ 5 observations )

Tree Pruning

→ A smaller tree with fewer splits ~~may~~ may lead to lower variance & better interpretation at cost of a little bias.

⇒ ~~This is~~ possible alternative build the tree only so long as the decrease in RSS due to each split exceeds some (high) threshold.

→ But this process is short-sighted since a seemingly worthless split early on in the tree might be followed by a very good split, i.e. a split that leads to a large reduction in RSS later on.

→ Better strategy is to grow a very large tree $T_0$, and then prune it back in order to obtain a subtree.

Goal: Subtree with lowest test error rate.

Method: Cost Complexity Pruning or (weakest link pruning)

↳ Algorithm 8.1 (ISLR 7th)

$\Rightarrow$ Cost function now is

$$\sum_{m=1}^{|T|} \sum_{x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha|T|$$

$|T| \Rightarrow$ indicates num. of terminal nodes.

$\Rightarrow$ Note:- For a higher value of $\alpha$, we get a tree which is sub tree of large tree $(\alpha = 0)$.

$\rightarrow$ use K-fold CV to find $\alpha$ that gives best test Error.

(FYI: Each value of $\alpha$ corresponds to one subtree)

② <u>Classification Trees</u>

$\rightarrow$ Assign most commonly occurring class of training observations in that region.

$\rightarrow$ (Classification Error) $= 1 - \max_{K} (\hat{P}_{mk})$

$\hat{P}_{mk}$ - proportion of training observations in the $m^{th}$ region that are from $k^{th}$ class.

$\rightarrow$ <u>Not sensitive to Tree-growing</u>, hence we use other alternative <u>Gini index</u>, <u>Entropy</u> etc.

$$= \sum_{K=1}^{K} \hat{P}_{mk} (1 - \hat{P}_{mk}) \rightarrow \text{measure of total variance across class}$$

If  **Relationship b/w  factors & Response**

is well approximated                Non-linear Relation
by a linear model                            ↓
            ↓                          Trees perform better
    linear Regression                  than Linear Model

$$(\theta = \delta)$$

get to → Popular Tree-based methods

     ⓐ → CART (Classification And Regression Trees)
              - uses Gini Index

     ⓑ → ID3 (Iterative Dichotomiser 3)
              - uses Entropy & Information Gain

     ⓒ → C4.5

get to → Interview questions on Decision Tree