# Classification Tree

Let us now look at various issues associated with construction of classification trees.

Let $T$ be a tree with $\tilde{T} = \{t_1, \ldots, t_M\}$ as the set of terminal nodes.

Let $\pi_{j(t)} \in \{\pi_1, \ldots, \pi_c\}$ denote one of the class labels, that is associated with terminal node $t$

A classification tree consists of $T, \tilde{T}, \{\pi_{j(t)}, t \in \tilde{T}\}$ and a partition $\{U(t) : t \in \tilde{T}\}$.

Let the learning sample be $\mathcal{L} = \{(x_i, y_i) : i = 1 (1) N\}$
$$(y_i \text{'s are the class labels})$$

$\underline{\forall t \in T}$, define

$N(t) = \#$ of sample patterns in $\mathcal{L} \ni x_i \in U(t)$

$N_j(t) = \# \quad - \quad - \quad - \quad - \quad - \quad - \quad - \ni x_i \in U(t)$
$$\text{and } y_i = \pi_j$$

Clearly, $\sum_j N_j(t) = N(t)$ & $\sum_{t \in \tilde{T}} N(t) = N$

$P(t) = \dfrac{N(t)}{N}$ ; estimate of $P(x \in U(t))$ based on $\mathcal{L}$

$P(\pi_j | t) = \dfrac{N_j(t)}{N(t)}$ ; estimate of $P(Y = \pi_j \mid x \in U(t))$

Recall that $\forall t \in T$



let $t_L = l(t)$ & $t_R = r(t)$

$$p_L^t = \frac{p(t_L)}{p(t)} \quad ; \text{ estimate of } P(\underset{\sim}{x} \in U(t_L) \mid \underset{\sim}{x} \in U(t))$$

$$p_R^t = \frac{p(t_R)}{p(t)} \quad ; \text{ estimate of } P(\underset{\sim}{x} \in U(t_R) \mid \underset{\sim}{x} \in U(t))$$

## Rule for label assignment (class label assignment)

Assign label $\pi_j$ to node $t$ if

$$p(\pi_j \mid t) = \max_i p(\pi_i \mid t)$$

(i.e. majority voting inside the partition)

## Splitting rules

Splitting rules at internal nodes are based on "node impurity" measures.

In general "node impurity" for any node $t \in T$ is

defined as

$$\text{Imp}(t) = \phi\left(p(\pi_1 \mid t), p(\pi_2 \mid t), \ldots, p(\pi_c \mid t)\right)$$

Note:

Such a "node impurity" would be maximum if

$p(\pi_j \mid t) = \frac{1}{c} \quad \forall j$ (i.e. equal representation of all classes at $t$)

and minimum if $p(\pi_j \mid t) = 1$ for some $j$ & $p(\pi_i \mid t) = 0 \quad \forall i \neq j$ — it's a pure node

Examples of "node impurity" measures

(i) Gini Index $(t) = \sum\limits_{\substack{i,j \\ i \neq j}} p(\pi_i | t) \, p(\pi_j | t)$

$\left(\begin{array}{l} \text{If } C = 2 \text{ (two-class problem)} \\[4pt] \text{Gini Index } (t) = \sum\limits_{i=1}^{2} \sum\limits_{\substack{j=1 \\ j \neq i}}^{2} p(\pi_i | t) \, p(\pi_j | t) \end{array}\right.$

$$= p(\pi_1 | t)\, p(\pi_2 | t) + p(\pi_2 | t)\, p(\pi_1 | t)$$

$$= p(1-p) + (1-p)\, p \qquad \left( p(\pi_1 | t) = p \right)$$

$$= 2\, p(1-p)$$

$\nearrow$ max if $p = \frac{1}{2}$ $\left( \text{i.e. } \frac{1}{c} \right)$

min if $p$ or $1-p = 0/1$ $\Big)$.

(ii) Misclassificate error rate at node $t$

$$= \frac{1}{N(t)} \sum\limits_{i:\, \underset{\sim}{x_i} \in U(t)} I\left( y_i \neq \pi_{j(t)} \right)$$

$j(t) = \arg\max\limits_{i} p(\pi_i | t)$

$$= \frac{1}{N(t)} \left( N(t) - N_{j(t)}(t) \right)$$

$$= 1 - \frac{N_{j(t)}(t)}{N(t)} = 1 - p(\pi_j | t)$$

(iii) Cross-entropy or deviance

$$- \sum\limits_{i} p(\pi_i | t) \, \log p(\pi_i | t)$$

<u>Remark</u>: The above measures are for node impurity.
We can define tree impurity as

$$\text{Imp}(T) = \sum_{t \in \tilde{T}} p(t)\,\text{Imp}(t)$$

## How to split a node ?

This is by far the most important question !!

Consider a split using variable $x_K$ at level $\ell$

say, $\mathcal{S}_K^\ell = \{\underset{\sim}{x} : x_K < \ell\}$

We define a measure of "goodness of split" at

node $t$ as the change in impurity $f^n$. For the

split $\mathcal{S}_K^\ell$ at $t$, this is

$$\Delta\,\text{Imp}(\mathcal{S}_K^\ell, t) = \text{Imp}(t) - \left( p_L^t\,\text{Imp}(t_L) + p_R^t\,\text{Imp}(t_R) \right)$$

$$\left( \text{recall that } p_L^t = P(\underset{\sim}{X} \in U(t_L) \,|\, \underset{\sim}{X} \in U(t)) \right.$$
$$\left. \& \; p_R^t = P(\underset{\sim}{X} \in U(t_R) \,|\, \underset{\sim}{X} \in U(t)) \right)$$

We need to find $(K^*, \ell^*) \ni \Delta\,\text{Imp}(\mathcal{S}_K^\ell, t)$ is

maximised over all $K = 1(1)p$ (dimension of feature vector)

and $\ell$ is allowed take one of a finite # of values

within the range of possible values of the

chosen feature variable.

The above "goodness of split" based approach is used to split nodes starting with the root node.

## When to stop splitting?

We do not split a node if the change in the impurity due to split is less than a prescribed threshold

### OR

Grow the tree till the terminal nodes are all pure (all patterns belonging to the partition have same class label) and then apply prunning of the tree.

## What is prunning? How to prune a classification tree?

There are various approaches of prunning. We discuss here 2 important approaches.

Prunning of a grown tree after applying splitting of nodes basically means cutting branches of the tree to get a subtree of the original tree.

## Cost-Complexity pruning

Let
$$r(t) = 1 - \max_i p(\pi_i | t).$$
↑
estimate of prob of misclassification at node $t$

define
$$R(t) = p(t) \, r(t)$$

Estimate of overall misclassification rate of the tree classifier is

$$R(T) = \sum_{t \in \tilde{T}} p(t) r(t) = \sum_{t \in \tilde{T}} R(t).$$

If $\alpha$ denotes the cost of complexity per terminal node then define

$$R_\alpha(t) = R(t) + \alpha \qquad \text{as cost-complexity}$$
$$\qquad\qquad\qquad\qquad\qquad \text{criterion for node } t$$

$$\& \quad R_\alpha(T) = \sum_{t \in \tilde{T}} (R(t) + \alpha) \quad \left( \begin{array}{l} \alpha : \text{tuning parameter} \\ \text{a trade-off parameter} \end{array} \right)$$

$$= \sum_{t \in \tilde{T}} R(t) + \alpha |\tilde{T}|$$

where $|\tilde{T}|$ : cardinality of terminal node set $\tilde{T}$.

Cost-complexity pruning approach : Starting from $T_0$ (a pure tree, i.e. a tree having all terminal nodes as pure), find the subtree $T_\alpha$ (for a fixed $\alpha$) $\ni R_\alpha(T)$ is minimum.

Weakest link pruning approach

Let $T_t$ be subtree with root $t$

& $\{t\}$ is subbranch of $T_t$ consisting of a single node $t$

Let $R_\alpha(t) = R(t) + \alpha = r(t)\, p(t) + \alpha$

(as defined earlier)

& for $T_t$; $R_\alpha(T_t) = R(T_t) + \alpha |\tilde{T}_t|$

Now $R_\alpha(T_t) < R_\alpha\{t\}$ $\left(\begin{array}{l} \text{i.e. } T_t \text{ has smaller} \\ \text{cost complexity than } \{t\} \end{array}\right)$.

if $R_\alpha(T_t) + \alpha |\tilde{T}_t| < R(t) + \alpha$

i.e. if $\alpha(|\tilde{T}_t| - 1) < R(t) - R(T_t)$

i.e. if $\alpha < \dfrac{R(t) - R(T_t)}{|\tilde{T}_t| - 1}$ $\qquad$ — (*)

Note that as $\alpha \uparrow$ equality is achieved at (*)

and $T_t$ and $\{t\}$ have same cost-complexity

and $\{t\}$ is preferred to $T_t$; i.e. $T_t$ can be

prunned at $t$

We define

$$g(t) = \dfrac{R(t) - R(T_t)}{|\tilde{T}_t| - 1}$$

as the "strength of link" at node $t$

Find $t^*$ $\ni$ $g(t^*)$ is min and then prune the tree at this "weakest-link", i.e. at $t^*$.

At the next step start with the prunned subtree and find the weakest link and prune again.

Continue to reach the root through prunning to get the sequence of prunned subtrees

$$T_1 \supset T_2 \supset \cdots \supset \{t_1\}$$

$\uparrow$ org root

Under this approach, the prunned subtree with min $R_\alpha(T)$ is the best prunned tree.

# Growing a regression tree

How to split nodes?

$$d = \{ (x_i, y_i) : i = 1(1)N \}$$

Consider a split variable $j$ and a split pt $t$

$$R_1(j,t) = \{ x : x_j < t \} \ \& \ R_2(j,t) = \{ x : x_j \geq t \}$$

Set the criterion $f^n$ as

$$\underset{j,t}{Min} \left[ \underset{c_1}{min} \sum_{i : x_i \in R_1(j,t)} (y_i - c_1)^2 + \underset{c_2}{min} \sum_{i : x_i \in R_2(j,t)} (y_i - c_2)^2 \right] - (*)$$

Note that for a fixed $j$ & $t$

$$min \to \hat{c}_1 = average \left( y_i \mid x_i \in R_1(j,t) \right);$$

$$\& \ \hat{c}_2 = average \left( y_i \mid x_i \in R_2(j,t) \right)$$

Find $(j^*, t^*)$ for the optimum split at a node $\ni$

$$\sum_{i : x_i \in R_1(j,t)} (y_i - \hat{c}_1)^2 + \sum_{i : x_i \in R_2(j,t)} (y_i - \hat{c}_2)^2 \quad \text{is minimised}$$

here $j$ varies over all $j = 1(1) p$ (feature vector dim)

& $t$ in the grid of the $j^{th}$ variable

## when to stop splitting?

(I) split tree node only if the decrease in sum of squares residual due to split exceeds some pre-assigned threshold

OR

(II) Std dev of responses at terminal nodes is low (below a threshold)

OR

(iii) Grow a large tree and stop splitting if the node size (# of patterns reaching that node) reaches a low threshold level and apply pruning.

## Regression tree pruning

$\forall t \in \tilde{T}$ define
$$\hat{c}_t = \frac{1}{N(t)} \sum_{x_i \in U(t)} y_i$$
or T

$N(t)$ : # of $x_i \in U(t)$

## Impurity measure:

Mean square error of obsn in $U(t)$ from $d$

$$Q_t = \frac{1}{N(t)} \sum_{i: x_i \in U(t)} (y_i - \hat{c}_t)^2$$

Define Cost complexity $f^n$ as

$$C_\alpha (T) = \sum_{t \in \tilde{T}} N(t)\, Q_t + \alpha |\tilde{T}|$$

$\alpha$: tuning parameter which governs the trade-off
bet^n tree size and goodness of fit

$|\tilde{T}|$: cardinality of $\tilde{T}$.

Prunning can be done based on $C_\alpha (T)$. For a
fixed $\alpha$, let $T_\alpha$ be a subtree of $T$ obtained

by prunning $T$. We try to find
$$T_\alpha \subset T \ni T_\alpha \text{ minimizes } \ell_\alpha (T)$$
for a fixed $\alpha$.

(Note: larger the $\alpha$ smaller in the opt $T_\alpha$)

<u>Weakest link prunning</u>: Collapse internal

nodes (to prune) that produces the smallest

per node increase in

$$\sum_{t \in \tilde{T}} N(t)\, Q(t) \text{ and continue till}$$

root node is reached.

Select the optimal from the prunned subtree
sequence.