

[illegible]

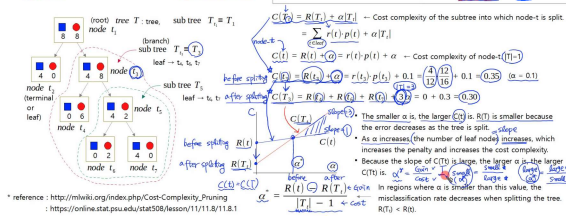
aims to find a subtree. Find a tree as good as the fully grown tree while having less # of leaves.

- Cost Complexity Pruning (CCP) : post-pruning – definition of Cost Complexity

- Since trying to lower the misclassification rate $R(T)$ results in a deeper tree, we define a new measure, $C(T)$ and build a tree that lowers $C(T)$.
 • $C(T)$ is called cost complexity. Objective: Min $C(T)$ instead of Min $R(T)$.
 • $C(T)$ is $R(T)$ plus penalty. The deeper the tree, the greater the penalty.
- misclassification rate of a tree $R(T)$ ↓ ↑ number of leaf node in the tree

• Cost Complexity: $C(T) = R(T) + \alpha T$

↓ ↑
 Penalty term



Minimal cost-complexity pruning is an algorithm used to prune a tree to avoid over-fitting, described in Chapter 3 of [BRE]. This algorithm is parameterized by $\alpha \geq 0$ known as the complexity parameter. The complexity parameter is used to define the cost-complexity measure, $R_\alpha(T)$ of a given tree T :

$$R_{\alpha}(T) = R(T) + \alpha|\bar{T}|$$

where $|T|$ is the number of terminal nodes in T and $R(T)$ is traditionally defined as the total misclassification rate of the terminal nodes. **Minimizing cost over all subtrees of T is equivalent to minimizing the cost of the terminal nodes.** As shown above, the cost of a node depends on the criterion. Minimal cost-complexity pruning finds the subtree of T that minimizes $R_{\alpha}(T)$.

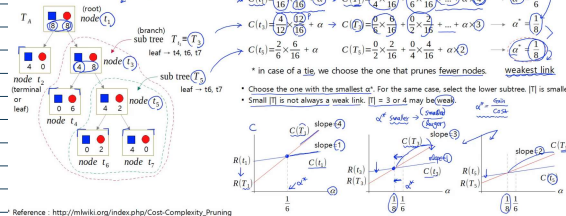
The cost complexity measure of a single node is $R_{\alpha}(t) = R(t) + \alpha$. The branch T_t is defined to be a tree where node t is its root. In general, the impurity of a node is greater than the sum of impurities of its terminal nodes, $R(T_t) < R(t)$. However, the cost complexity measure of a node, t and its branch, T_t can depend on each other. We define the effective cost of a node to be the value where they are equal, $R_{\alpha}(t) = R_{\alpha}(T_t) = R(t) - \alpha$ or $\alpha_{eff}(t) = \frac{R(t) - R(T_t)}{2}$. A non-terminal node with the smallest value of α_{eff} is the weakest link and will be pruned. This process stops when the pruned tree's minimal α_{eff} is greater than the `ccp_alpha` parameter.

- Cost Complexity Pruning (CCP) : post-pruning – Example

- **Step-1:** Build a tree as deep as possible and find the subtree with the smallest α^*
- tree T : tree, sub tree $T_i = (T_i)$
- $$C(t) = R(t) + \alpha$$

$$C(T_i) = R(T_i) + \alpha |T_i|$$

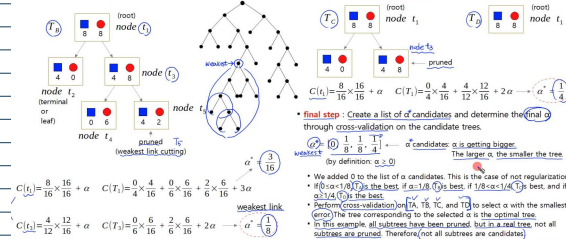
$$\alpha^* = \frac{R(t) - R(T_i)}{|T_i| - 1}$$

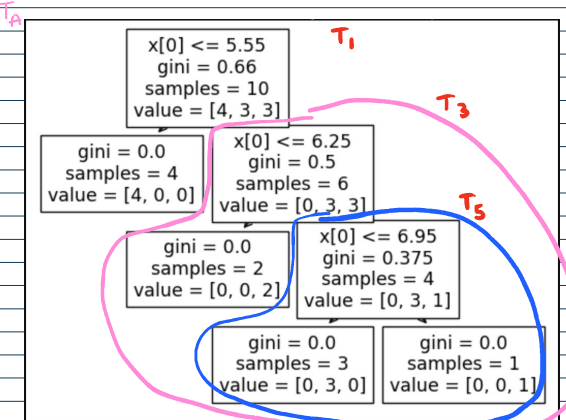


$$\Rightarrow \alpha^* = \frac{R(H) - R(L_n)}{1 + 1} = \frac{\text{Decrease in Impurity caused by split}}{\# \text{ of leafs added to the tree}} = \frac{\text{gain}}{\text{cost}}$$

- Cost Complexity Pruning (CCP) : post-pruning – Example

- Step-2 :** Prune the sub tree found in step-1, and find the sub tree with the smallest α again.
- Step-3 :** Prune the sub tree found in step-2, and find the sub tree with the smallest α . Repeat this process.





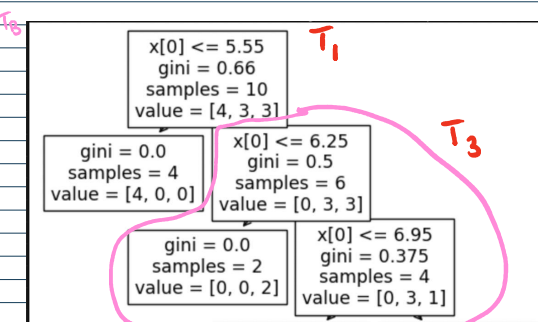
$$R(t_1) = \frac{10}{10} \times 0.46 + \alpha \quad R(T_1) = 0.4\alpha \quad \alpha^* = \frac{R(t_1) - R(T_1)}{4-1} = 0.22$$

$$R(T_3) = \frac{6}{10} \times 0.5 + 1 \quad R(T_3) = 0.3 + 1 \quad \alpha^* = \frac{R(T_3) - R(T_2)}{3-1} = 0.15$$

$$R(t_5) = \frac{4}{10} \times 0.275 = 0.11 \quad R(T_5) = 0.275 \quad \alpha^* = \frac{R(T_5) - R(t_5)}{2-1} = 0.15$$

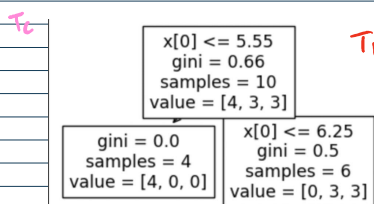
$$N_t / N * (\text{impurity} - N_{t_R} / N_t * \text{right_impurity} - N_{t_L} / N_t * \text{left_impurity})$$

↳ How much decrease in impurity this particular splitting has, compared to the whole tree.

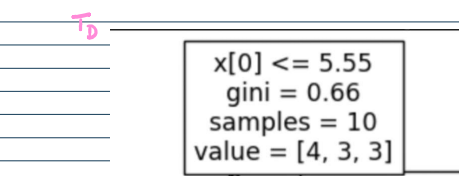


$$R(T_1) = \frac{10}{10} \times 0.45 + \alpha \quad R(T_2) = \frac{4}{10} \times 0.325 + 3\alpha \quad \alpha = \frac{R(T_1) - R(T_2)}{3-1} = 0.255$$

$$R(T_2) = \frac{6 \times 0.5 + 2}{10} \quad R(T_2) = \frac{4 \times 0.5 + 2}{10} \quad \alpha^T = \frac{R(T_2) - R(T_1)}{2 - 1} = 0.15$$



$$R(t_1) = \frac{10}{10} \times 0.66 + \alpha \quad R(t_2) = \frac{6}{10} \times 0.5 + 2\alpha \quad \alpha = \frac{R(t_2) - R(t_1)}{2-1} = 0.36$$



$$\alpha^* = \begin{matrix} & T_A & T_B & T_C & \rightarrow T_D \\ \left[\begin{matrix} 0 & 0.15 & 0.15 & 0.36 \end{matrix} \right] \end{matrix}$$

in sklearn $\rightarrow \alpha^* = [0, 0.15, 0.36]$

T_A T_C T_D

```
path['ccp_alphas']  
array([0. , 0.15, 0.36])  
  
path['impurities']  
array([0. , 0.3 , 0.66])
```

$T(\text{ccp_alpha} = X)$



- ① start from fully grown tree
- ② prune the tree w.r.t minimum α
- ③ stop when pruned tree has minimum $\alpha > X$

$G.I(T)$ = weighted sum of
 $G.I$ leaf nodes