

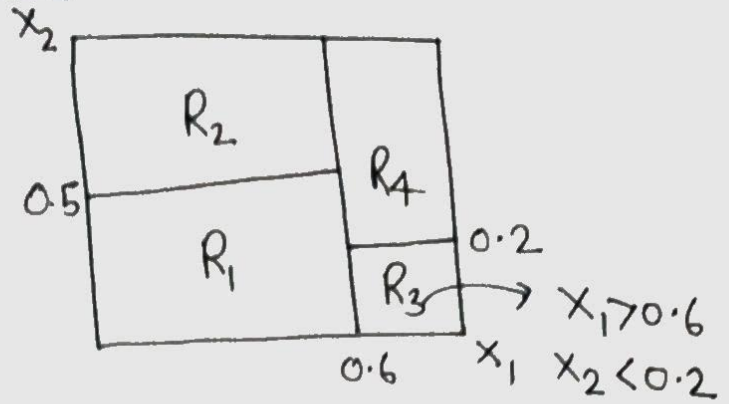
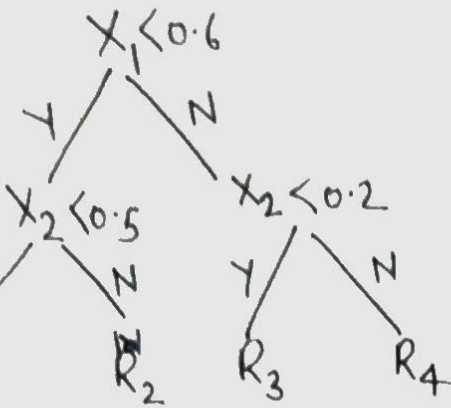
Lecture-33-1 Decision Trees - Introduction

Querying properties of the Data

— Partition I/P space into rectangles.

In theory regions could be of any shape

we use high dimensional rectangular boxes.



- highly Interpretable
- Universal approximator
- Non-parametric

Regression Trees

— Fit a constant in each region.

$$\hat{f}(x) = \sum_{m=1}^4 c_m I\{(x_1, x_2) \in R_m\}$$

$(x_i, y_i), i=1, \dots, N, x_i \in \mathbb{R}^p, y_i$

$x_i = \{x_{i1}, x_{i2}, \dots, x_{ip}\}$

$$\Rightarrow f(x) = \sum_{m=1}^M c_m I(x \in R_m)$$

$$= c_1 I(x \in R_1) + c_2 I(x \in R_2) + c_3 I(x \in R_3) + c_4 I(x \in R_4)$$

— Determine the 'M' regions

— Given regions, find Response. $c_m \rightarrow$ response for region 'm'.

Minimise RSS:

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

$$\hat{y}_{R_m} = \hat{c}_m = \frac{1}{|R_m|} \sum_{i \in R_m} f(x_i) \rightarrow \text{Mean response For training observations}$$

Lecture-34-1

Approach (greedy)

- $$R_1(j, s) = \{x | x_j \leq s\}, \quad R_2(j, s) = \{x | x_j > s\}$$

$$\min_{J, S} \left[\sum_{i: x_i \in R_1(J, S)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2(J, S)} (y_i - \hat{y}_{R_2})^2 \right]$$

\downarrow $C_1 \rightarrow$ Mean response
 \downarrow $C_2 \rightarrow$ Mean response

we will do separately.

$$\underline{j \in \{1, 2, \dots, p\}}$$

fix J and find S

Choose 'S'.

$n \rightarrow$ Total points. $\} \text{ Total time: } \boxed{n * p} \rightarrow \text{For one split}$
 $p \in J$

↓ → Can be repeated while growing trees.

Lecture 35 - Stopping criteria & pruning

*Early stopping - May lead into poor decision boundary.

*Leaf of a tree is a region

Tree pruning

Above process will lead to over fitting

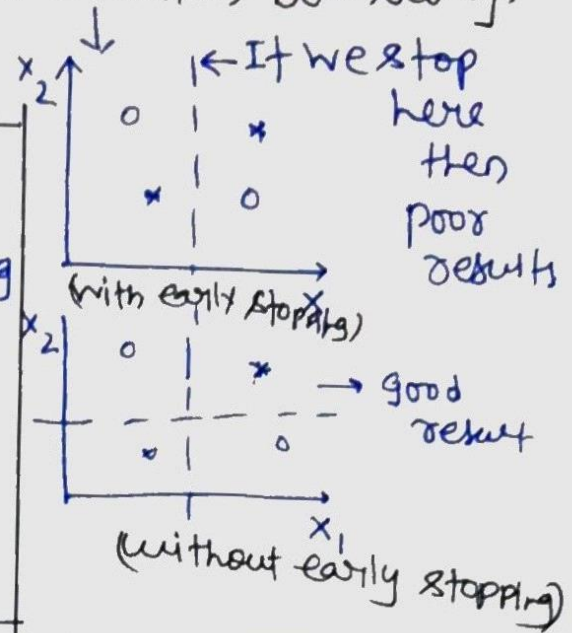
→ Tree will be very complex.

→ Variance higher

→ Simple Tree

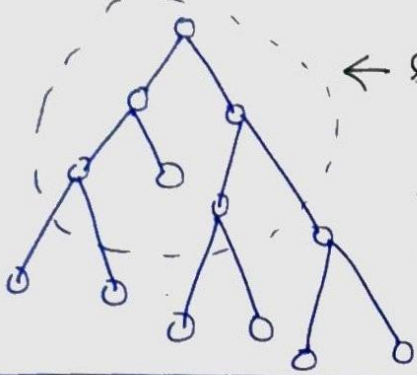
less regions → lower variance.

better interpretation



Strategy: Grow Tree T , and then prune it back in order to obtain subtree. using (cross-validation)

← selecting best subtree will be very time consuming → No. of subtrees are very large.



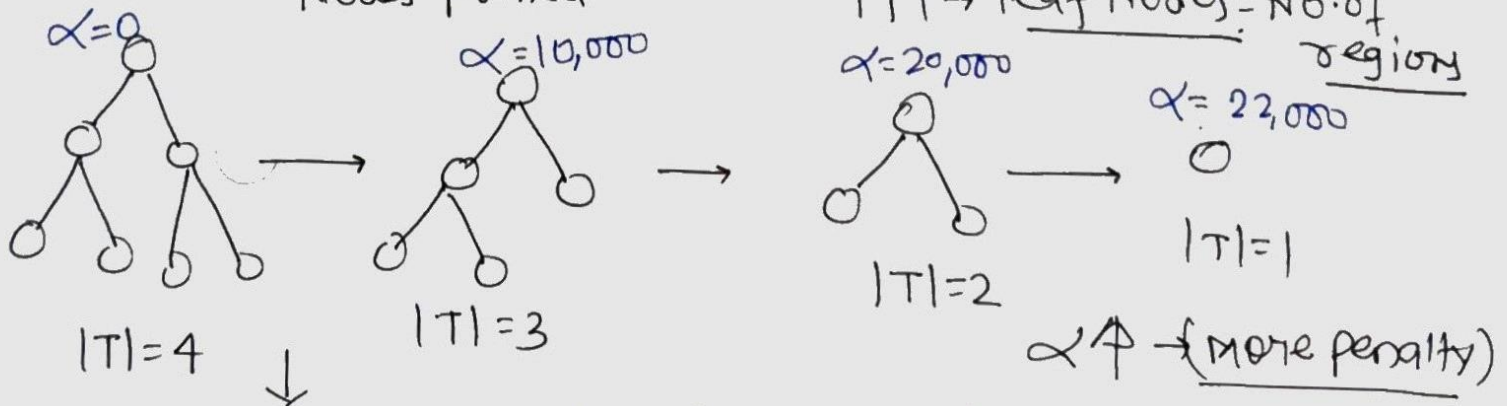
Cost-Complexity pruning: weaker link pruning

$$C_\alpha(T) = \text{prediction error} + \alpha |T|$$

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

Complexity penalty

As $\alpha \uparrow \Rightarrow$ Number of Nodes pruned.



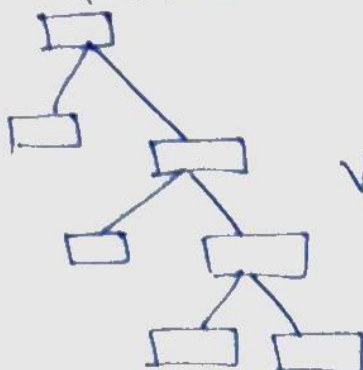
~~We want to penalise tree more having less no. of regions~~

~~To Compare~~

① we already know that pruned tree will have more errors.

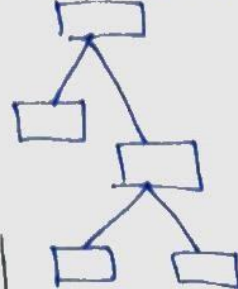
② we need to introduce one more term by which we can compare these trees. by introducing Complexity penalty.

$$SSR = 543.8$$



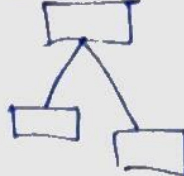
VS

$$SSR = 5494.8$$



VS

$$SSR = 19243$$



VS

$$SSR = 28897.2$$



Let's take $\alpha = 10,000$

$$\text{Total Score} = SSR + 10,000 \cdot T$$

$$\text{Tree Score} = 40,543.8$$

$$\text{Tree Score} = 35,494$$

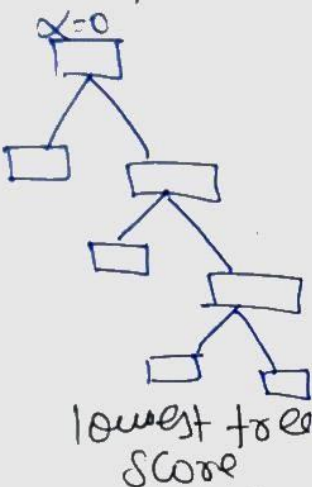
$$\text{Tree Score} = 39,243.1$$

$$\text{Tree Score} = 38,897.2$$

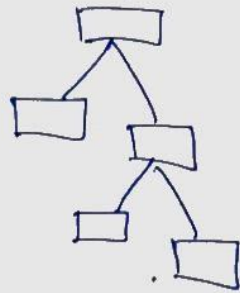
Optimal Subtree

Using different α can give different result.
We need to find best value for alpha.

$\alpha = 0$, \Rightarrow Full sized tree will have lower Tree Score.

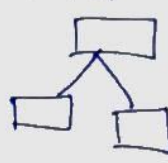


$\alpha = 10,000$



lower tree score

$\alpha = 15,000$



(lower tree score)

$\alpha = 22,000$



\downarrow
gives lower tree score.

Use testing data to get the optimal subtrees for different values of $|T|$.

Lecture 36

Decision Tree for classification - Loss function.

$$\hat{p}_{mk} = \frac{1}{|R_m|} \sum_{x_i \in R_m} I\{y_i \neq k\}$$

Total Number of Correct predictions

↓
Total Number of data points belonging to class k in region m

↓
prediction

2 classes	
R_2	R_3
R_1	R_4

$p_{41}, p_{42}, p_{43}, p_{44}, p_{21}, p_{22}, p_{23}, p_{24}$
 $p_{11}, p_{12}, p_{13}, p_{14}, p_{31}, p_{32}, p_{33}, p_{34}$

if $x_i \in R_m$ then check which is greater.
 $[p_{41}, p_{42}, p_{43}, p_{44}]$

$$K(m) = \arg \max_K \hat{p}_{mk} = \max_K [p_{41}, p_{42}, p_{43}, p_{44}]$$

$$K(4) =$$

↓ label for region 4

$$\text{classification error} = \frac{1}{|R_m|} \sum_{x_i \in R_m} I(y_i \neq K(m)) = (1 - \hat{p}_{m(K(m))})$$

$$= (1 - \max_K \hat{p}_{mk})$$

Gini Index := $\sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk})$ - Measure of variance across K -classes for a particular region.

Cross entropy: $-\sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$ (information gain criteria)

↓
 estimated label

Lecture-37.

Decision Trees - Categorical attributes

q- unordered attributes:

Divide into 2 groups. Total No. of grps possible:

$$(q_1 + q_2 + q_3 + \dots + q_{q/2}) = 2^{q-1} - 1 \rightarrow \text{Not feasible to go and pick split points.}$$

in case of regression: you have to look 'n' points.
n → No. of data points.

0/1 - binary classification.

Suppose $q=5$ For color attribute:

$q=0$ Red,

$q=1$ green

$q=3$ blue

$q=4$ yellow

$q=5$ magenta.

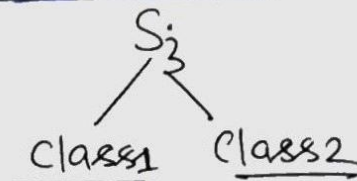
$P(\text{class1}|\text{Red}), P(\text{class1}|\text{green}),$

$P(\text{class1}|\text{blue}), P(\text{class1}|\text{yellow}),$

$P(\text{class1}|\text{magenta}) \rightarrow$ Ascending order. and then make the split.

0.2 | 0.3 | 0.4 | 0.45 | 0.55
R | G | B | Y | M

$S_1 \quad S_2 \quad S_3 \quad S_4$ ← split point.

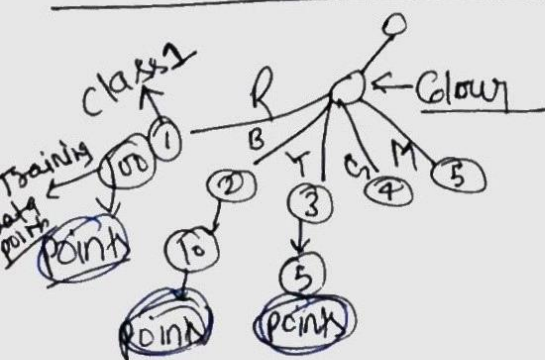


No-Need to consider $(2^{q-1} - 1)$ splits.

Multi-class

Heuristics: -

Multiway-splits



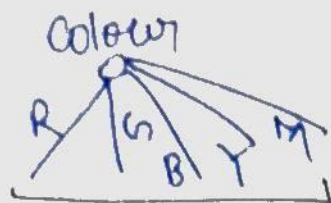
Lecture-38. Multiway splits.

Problem

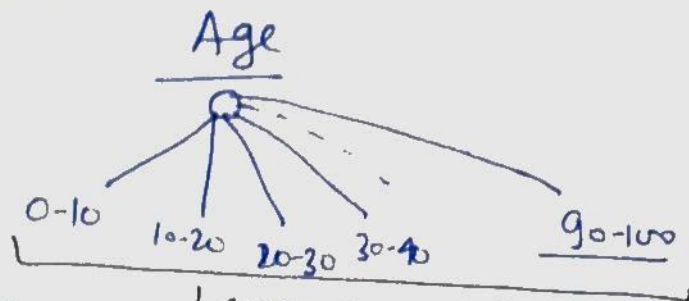
- spawling (spreading) (hard to interpret)
- lost interpretability
- Variance max.
- At each split, data points will be lower..
- Tree will become sparse.
- Overfitting is the problem.

Lecture-38

— Favours attributes with more values.



$\text{Error}_{\text{Colour}} > \text{Error}_{\text{Age}}$



Less error in splitting

C4.5 - Algorithm - gain-Ratio

↓
What is the information gain
after splitting into 10 parts
rather than 6.

→ we can also use binary split to arrive at same
→ function which is predicted by multiway splits.

binary split is always better > multiway
split ~~because~~
↓
we can use recursive approach to make
multi-way split.

Multiway splits → Avoids the choosing variable (5)
(split point)