

DECISION

classmate

Date _____

Page _____

TREE

(CLASSIFICATION
& REGRESSION)

↳ NESTED IF-ELSE.

- Pure Node : Node in Decision Tree where all datapoints belong to the same class. or have same value (regression)
→ leaf node are Pure Node
- Impure Node : where all datapoints don't belong to same class or same value

AIM OF DT : Start of Impure node to reach to Pure Node

★ In some cases, it can happen we can't reach to pure state

Q How do we calculate to take which feature next in order to reach pure state? INFORMATION GAIN

→ we use ~~loss~~ ~~for~~ ~~or~~ ~~test~~ ~~Pr~~ ~~n~~

Loss Function \Rightarrow i) ENTROPY
 in D.T ii) GINI INDEX

LOSS FUNCTION

① Entropy (0 to 1)

$$H(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

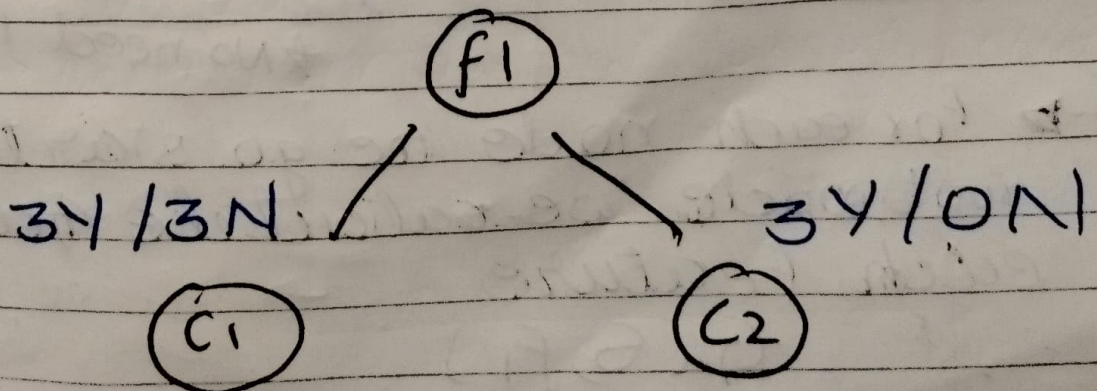
$p_+ \Rightarrow$ probability of Yes

$p_- \Rightarrow$ probability of No.

★ Pure Split \Rightarrow Entropy $\Rightarrow 0$

Impure Split \Rightarrow Entropy $\Rightarrow 1$

② Check which node is Pure and which is Impure? (C_1 & C_2)



For C_1

$$H(S) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$H(S) = 1$$

↳ Impure split

For C_2

$$H(S) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$

$$H(S) = 0$$

↳ Pure split.

FEATURE TO TAKE NEXTINFORMATION GAIN

$$\text{Gain}(S, f_x) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v)$$

↓
feature
v ∈ val
|S_v|

* No need to understand formula

---> For each node we go starting from root node we calculate Gain for each feature

$$G(S, f_1)$$

$$G(S, f_2) \dots$$

then select that feature there which has the highest Gain.

LOSS FUNCTION \rightarrow GINI INDEX (DEFAULT)

$$G.I = 1 - \sum_{i=1}^n (p_i)^2$$

$$p^2 = (p_+)^2 + (p_-)^2$$

$$\bigcirc \quad 2Y/2N$$

$$G.I = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = \frac{1}{2}$$

$G.I \Rightarrow 0.5 \Rightarrow$ Impure Split

$G.I \Rightarrow 0 \Rightarrow$ Pure Split.

Q In D.T which loss fxn to use when?
 \rightarrow When you have a large D.T with many parameter use Gini Index as loss fxn as it has less T.C as compare to Entropy which has log fxn in it \rightarrow

- ★ Decision Tree are very good at which column to include or not

Procedure

→ Based on Gain

1. Choose the Best Feature
2. Split the dataset of that category
- y in category.
3. Repeat.
4. Stop when the leaf node is pure.

Q How to select the Best Feature?
at each node?

↳ Two way

(ID3)

↳ If you using Entropy
as loss fn

↳ Calculate I.G for
each feature & select
the one with highest
I.G

(CART
Algo)

↳ G.I (Default)

↳ Select feature
with Gini Index

DECISION TREE

REGRESSOR

↳ Loss Fxn \Rightarrow Mean Square Error.

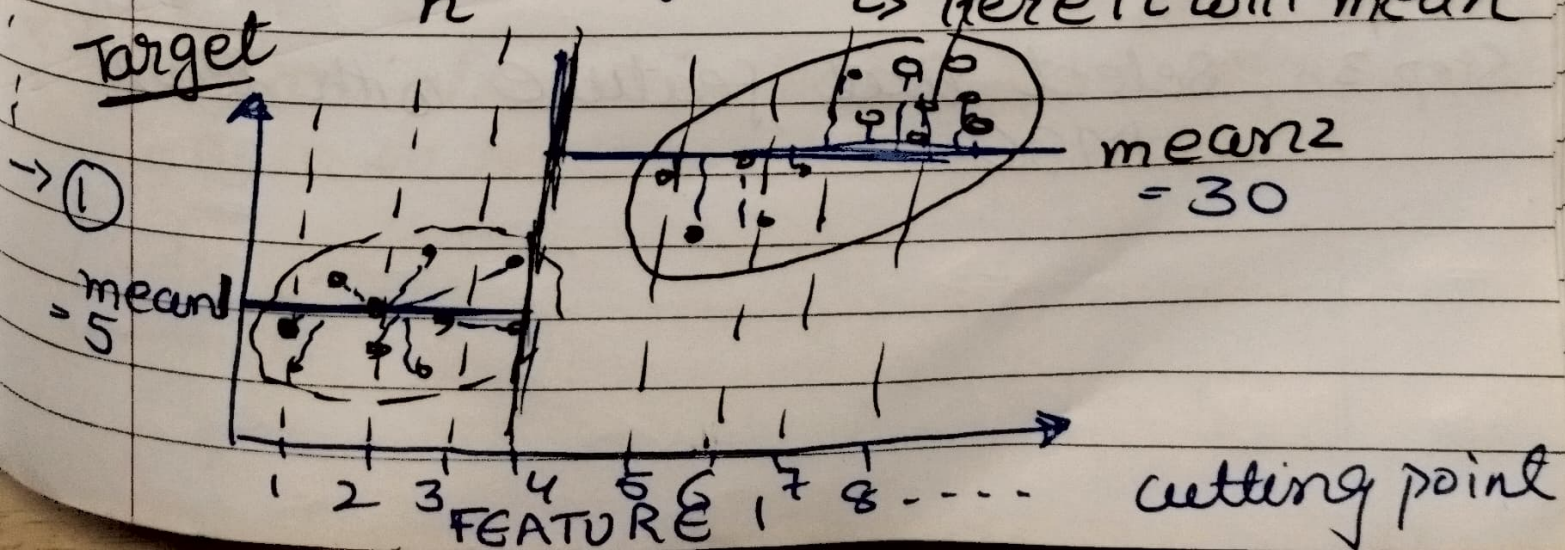
Q How we select the best feature in Regression?

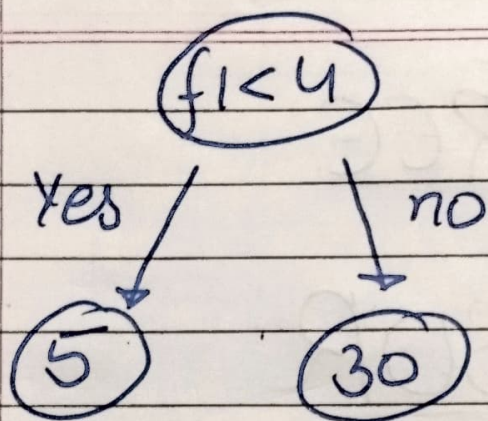
- ① First we calculate least MSE of each feature
- ② Then select that feature among them which have the least MSE

MSE

$$\Rightarrow \frac{1}{n} \sum (y - \text{pred})^2$$

↳ Here it will mean





$$\begin{aligned}
 \text{Error} &= \text{MSE1} + \text{MSE2} \\
 &\quad \text{(mean1)} \quad \text{(mean2)} \\
 &\quad f1 < 4 \quad f > 4
 \end{aligned}$$

---> Error

Step 1 Calculating Least MSE for a Feature

- i) Calculate MSE for each cutting point
- ii) Select the cutting point with least MSE.

Step 2: Do This for all feature

Step 3: Select that feature with least MSE.



Problem with D.T

↳ ① Overfitting

↳ i) Tree Pruning

ii) Use Random Forest.
(Ensemble Method)

★ `from sklearn.tree import DecisionTreeClassifier`

★ → To See Graph (D.T)

↳ `from sklearn import tree`
`plt.figure(figsize=(15, 10))`
`tree.plot_tree(model, filled=True)`

★ You can also check feature importance in Decision Tree

Code