Decision Tree Handle Theory
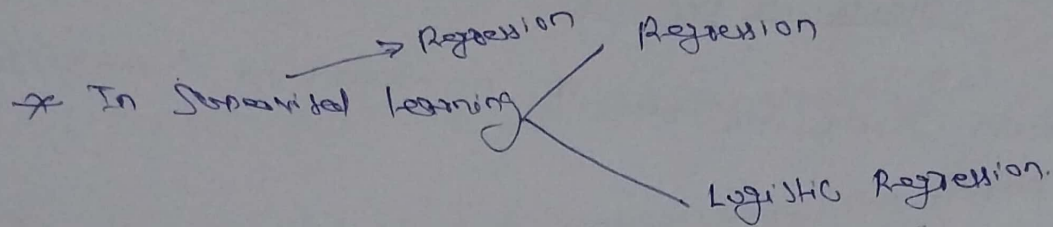
\* In Supervised learning → Regression, Regression

Logistic Regression.

\* if output variable is yes or not mean two only
It is Call binary Classification.

→ if output is more than two its called multi classi-
fication.

→ The benefit of Decision Tree is it Can handle
both numerical and Categorical data of.

for Example

Age → Numerical
Salary → Numerical & Continuous value.
Gender → Classification.
Occupation → Classification

→ Note !—
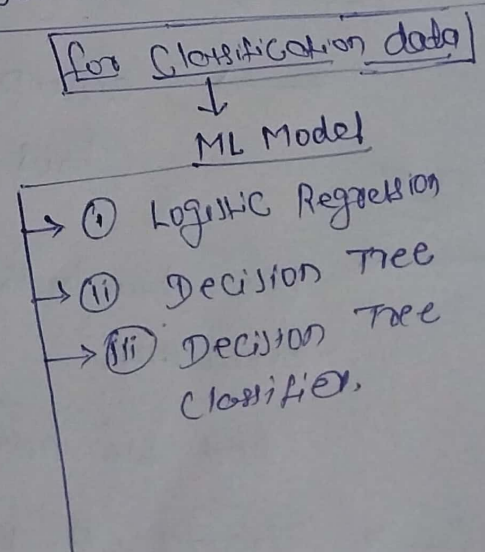Output here is binary Classification.

formula of entropy:—

$$= \sum_{j=1}^{2} Pi \log_2 Pi$$

how to Calculate entropy values.

$$- [0.5 \log_2 0.5 + 0.5 \log_2 0.5]$$

→ In Machine Learning entropy
exist as a measures of uncertainty.

for Classification data
↓
ML Model

→ ① Logistic Regression
→ ① Decision Tree
→ ①i Decision Tree
Classifier.

## Calculation Entropy

→ Decision Tree works on the principal or homogeneouness or purity.

→ Decision tree overcome the limitation of Logistic Regression,
    * It, is create only one decision boundary.
    So Decision trees come into the picture.
    there are several data which need many decision boundary.

→ Graph of data is for Algorithm, and tree is form.

→ ID3 ( Iterative Dichomizer 3) → old
    → best for Classification datas.

→ if the less entrophy is there means in data more homogenity is there

→ if high then less homogenity is in the data

→ entropy is the measure of Randomness or Chaos, mixture.

53.54

→ One drawbacks of decision tree is it 90% Stuck with overfitting.
    And we have to Control it.

→ Scarcity is very useful in case of inbalanced data.

→ Sectrons/ Box / Decision boundary

→ Over learning is also known as overfitting.

→ We should always ensure that our training accuracy never ever be 100%.

→ Whenever extra node added into the tree the it will be always bottom of tree

───────────→ Confusion Matrix

**PREDICTED**

| 165 | NO | Yes. |
|-----|------|------|
| NO | TN (50) | FP (10) |
| Yes | FN (5) | TP (100) |

(left side label: A C T U A L)

Accuracy Score $= \dfrac{TP + TN}{Total}$

$= \dfrac{50+100}{165}$

$= 0.91$

Error Rate $= 1 - accuracy$

$= 1 - 0.91$

$= 0.09$

or $= \dfrac{FP + FN}{TN+FP+FN+TP}$ (Total)

$= \dfrac{10+5}{165}$

$= 0.09$

Recall :- $\dfrac{TP}{actual\ Yes}$

$= \dfrac{100}{105}$

$= 0.95$

precision $= \dfrac{TP}{predicted\ Yes}$

$= \dfrac{100}{FP + TP}$

$= \dfrac{100}{10+100}$

$= \dfrac{100}{110}$

$= 0.64$

To develop such a model, the computed information gain (C, pitch) with respect to target is ---- (rounded off to two decimal places).

| Match No. | pitch | Format | Winner (Target) |
|-----------|-------|--------|-----------------|
| 1 | S | T | Green |
| 2 | S | T | Blue |
| 3 | F | O | Blue |
| 4 | S | O | Blue Blue |
| 5 | F | T | Blue Green |
| 6 | F | O | Green Blue |
| 7 | S | O | Blue Green |
| 8 | F | T | Green Blue |
| 9 | F | O | Blue Blue |
| 10 | S | O | Blue Green. |
|    |   |   | G+ |

Entropy Calculation

$$Entropy(S) = -P(yes) \log_2 P(yes) - P(no) \log_2 (P_{no})$$

where;

— S is the total sample space
P(yes) is probability of yes.

if number of yes = number of no ie P(S) = 0.5

⟹ Entropy(S) = 1

if it contains all yes or all no ie (P(S) = 1 or 0

⟹ Entropy(S) = 0

$$E(S) = -p(ye) \log_2 p(yes) - p(no) \log_2 (p(no))$$

$$= -0.5 \log_2 0.5 - 0.5 \log_2 0.5$$

$$= -0.5 (\log_2 0.5 - \log_2 0.5)$$
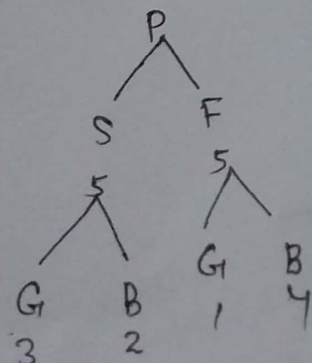
$$=$$

$$= 1$$

→

- Step 1: Entropy of entire datasets

$$G\{4,6\} = -\frac{4}{10} \log_2 \frac{4}{10} - \frac{6}{10} \log_2 \frac{6}{10} = 0.97$$

Step 2: Entropy of all attributes in pitch

$$\{3,2\} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$$

$$IG(p) = entropy \ (whole)$$

Entropy of Spin (S)

STEP 2

Progres SQL | Microsoft our training, Day of best
free. → Check now

Entropy of pace (F)

$$\{1,4\} = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} = 0.72$$

Information Gain (C, Pitch) = $Entropy(G) - \frac{5}{10} Ent(S) - \frac{5}{10} Ent(F)$

$$= 0.97 - 0.485 - 0.360$$

$$= 0.13$$

* what is gini_Index

Gini idex measures the impurity of the data.
highon Gini meos more impurity.

11.38.09