① Supervised Learning: Supervised learning is a type of machine learning where the model is train using labeled data — that means each input has a corresponding correct output (label).

labeled data example

| customer-id | Age | gender | Income (Lakh) | Years-employed (years) | total wealth Lakh |
|---|---|---|---|---|---|
| 1021210 | 27 | Female | 1 lakh | 1 | 10 |
| 20306212 | 28 | male | 2 lakh | 10 | 20 |

Types of algorithm used in supervised learning:

classification

① Logistic Regression

② K-Nearest Neighbour (KNN)

③ Decision Tree Classifier

④ Random forest classifier

⑤ Naive Bayes classifier

⑥ SVM (support vector machine

⑦ Gradient Boosting (XGBoost, LightGBM, CatBoost)

⑧ NN (ANN, CNN, RNN)

Regression

① Linear regression

② Ridge / Lasso Regression

③ Polynomial regression

④ Decision tree Regression

⑤ Random forest u

⑥ SVR ⑦ XGBoost, LightGBM

**Unsupervised learning** → Unsupervised learning is a type of machine learning to where model is trained using underlabeled data - meaning there is no pre-define output.

Examples :

(i) Grouping customer behavior. (Clustering)
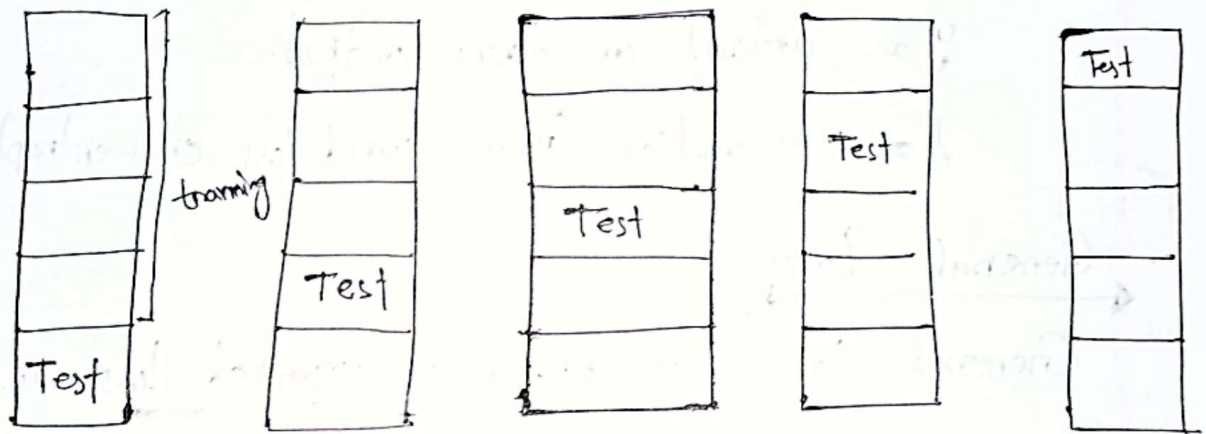
(ii) Reducing dimentions for visualisation (PCA)

Some examples of unsupervised learning algorithm

(i) K-mean clustering

(ii) Hierarchical clustering

(iii) DBSCAN (Density-Based Spatial clustering)

(iv) GMM (Gaussian Mixture Models)

K- Fold cross validation: → K- fold Cross Validation is a model evaluation tahnique used to how well a machine learning model perform on unsean data.

It works dividing the dataset into K equal folds (parts) then training and testing the model K-times, each time using a different fold as the test set and the remaining folds as the training set. Let's assume, K=5 then



Dataset
↓

Accuracy   $a_1$         $a_2$         $a_3$         $a_4$         $a_5$

K-fold - cross validation score = $\dfrac{a_1 + a_2 + a_3 + a_4 + a_5}{5}$

Empirical loss · also called · training loss :

Empirical loss is the average loss of a model over the # training dataset. It measure how well the model fits the training data.

$$L_{empirical} = \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i)$$

$n$ = number of training sample / batch size

$f(x_i)$ = model prediction for input $x_i$

$y$ = actual or true output.

$\ell$ = loss function (mean-squared loss, cross-entrophy loss).

General Loss:

General loss also known as expected loss on generalization Loss. General loss is the expected value of the loss over the entire data distribution, including unseen and future data. It measure how well a model generalizes or perform to new, unseen dataset.
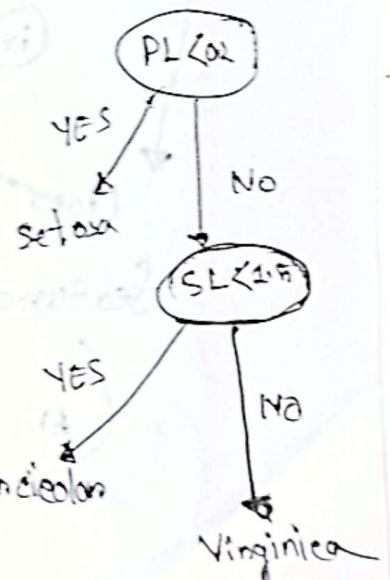
$$L_{general} = \mathbb{E}_{(x,y) \sim P_{data}} [\ell(f(x), y)]$$

## Difference

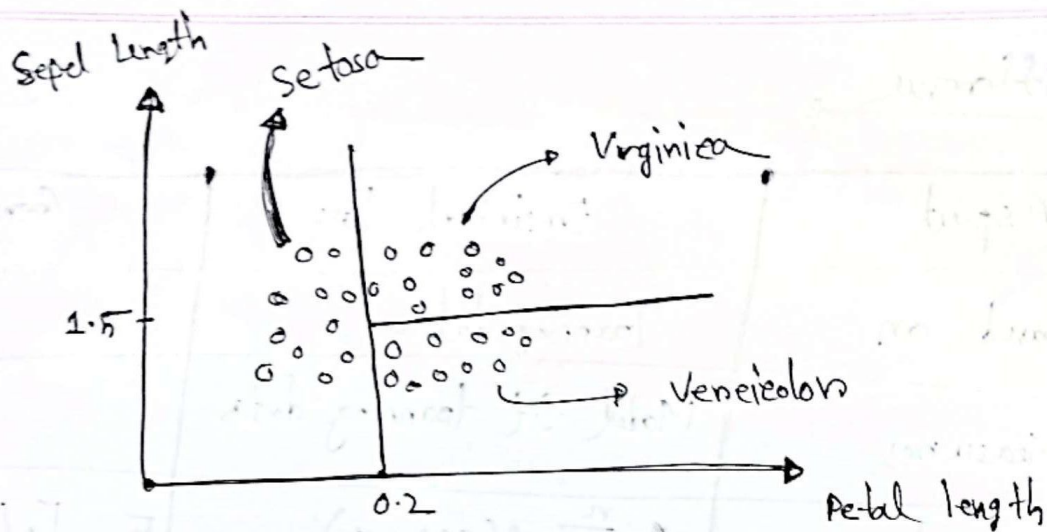| Aspect | Empirical Loss | General (Expected) Loss |
|---|---|---|
| Based on | Training data | |
| Measures | Model fit training data | |
| Formulas | $\frac{1}{n}\sum_{i=1}^{n} l(f(x_i), y_i)$ | $E_{(x,y)}\left[l(f(x), y)\right]$ |
| Depends on | The specific training data | Real-world data distribution |
| Risk | Low empirical loss may cause overfitting | Low general loss means good generalization |

## Decision Tree (classification Problem)

| Petal Length (PL) | Sepal Length (SL) | Type |
|---|---|---|
| 1·34 | 0·34 | Setosa |
| 3·45 | 1·45 | Versicolor |
| 1·62 | 0·98 | setosa |
| 2·56 | 1·79 | Virginica |
| 3·00 | 1·13 | versicolor |
| 1·3 | 0·88 | Setosa |



[ Decision tree model is nothing
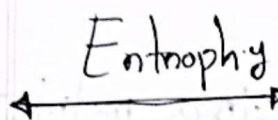more than a collection of
nested if-else statement. ]

Decision trees steps

(i) Dataset

(ii) Best features (Previous problem 1st PL than SL)

(iii) Split data base on best features

(iv) Repeat (1, 2, 3)
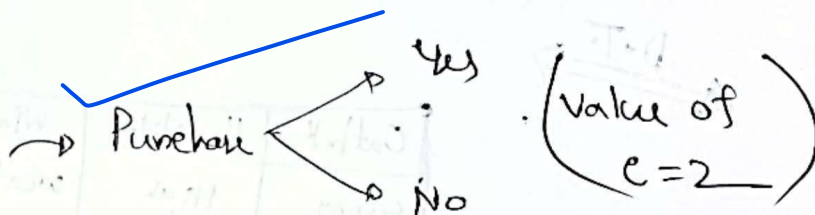
Now our question is how to find the best features?

ANS [ with entrophy and information gain.]

← Entrophy → Entrophy is nothing but the measure of disorderness or the measure of impurity. The mathematical formula of entrophy is

$$E(s) = \sum_{i=1}^{c} -P_i \log_2 P_i$$

$P_i$ is the frequency probability of an element/class is in our data.

| Salary | Age | Purchase |
|--------|-----|----------|
| 20000 | 21 | Yes |
| 10000 | 45 | No |
| 60000 | 27 | Yes |
| 15000 | 31 | No |
| 12000 | 13 | No |

→ Purchase < Yes ... No   (value of $c=2$)

$$E(d) = \sum_{i=1}^{2} -P_i \log_2 P_i$$

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

$$= -\left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) - \left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right)$$

$$= 0.97$$

## Information Gain :

Information gain is a metric used to → Train Decision Trees. Information gain is used in decision trees to find the best attributes/column to split the data at each node. Formula

* Information Gain = Entrophy — Weighted Entrophy

Parenty on target column entrophy

কারণ ও বায়

### D.T.

| Outlook | Humidity | Wind | Play tennis |
|---------|----------|--------|-------------|
| sunny | High | weak | No |
| sunny | High | strong | No |
| Rain | High | strong | No |
| Rain | Normal | weak | Yes |
| Rain | Normal | strong | No |
| sunny | Normal | strong | Yes |

**Step: 1**

entropy of Playtennis $= -P_{no} \log_2 P_{no} - P_{yes} \log_2 P_{yes}$

$$= -\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right) = 0.92$$

**1. outlook column**

$$E_{sunny} = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) = 0.92$$

$$E_{Rain} = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.92$$

$$\text{Weighted, } E = \left(\frac{3}{6} \times 0.92\right) + \left(\frac{3}{6} \times 0.92\right) = 0.92$$

$$IG_{outlook} = 0.92 - 0.92 = 0$$

**2. Humidity**

$$E_{High} = -\frac{3}{3} \log_2\left(\frac{3}{3}\right) - \frac{0}{3} \log_2\left(\frac{0}{3}\right) = 0$$

$$E_{normal} = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log\left(\frac{1}{3}\right) = 0.92$$

$$\text{IG Humidity}$$

$$\text{Weighted, } E = \left(\frac{3}{6} \times 0\right) + \left(\frac{3}{6} \times 0.92\right) = 0.459$$

$$IG_{Humidity} = 0.92 - 0.459 = 0.461$$

**3. wind**

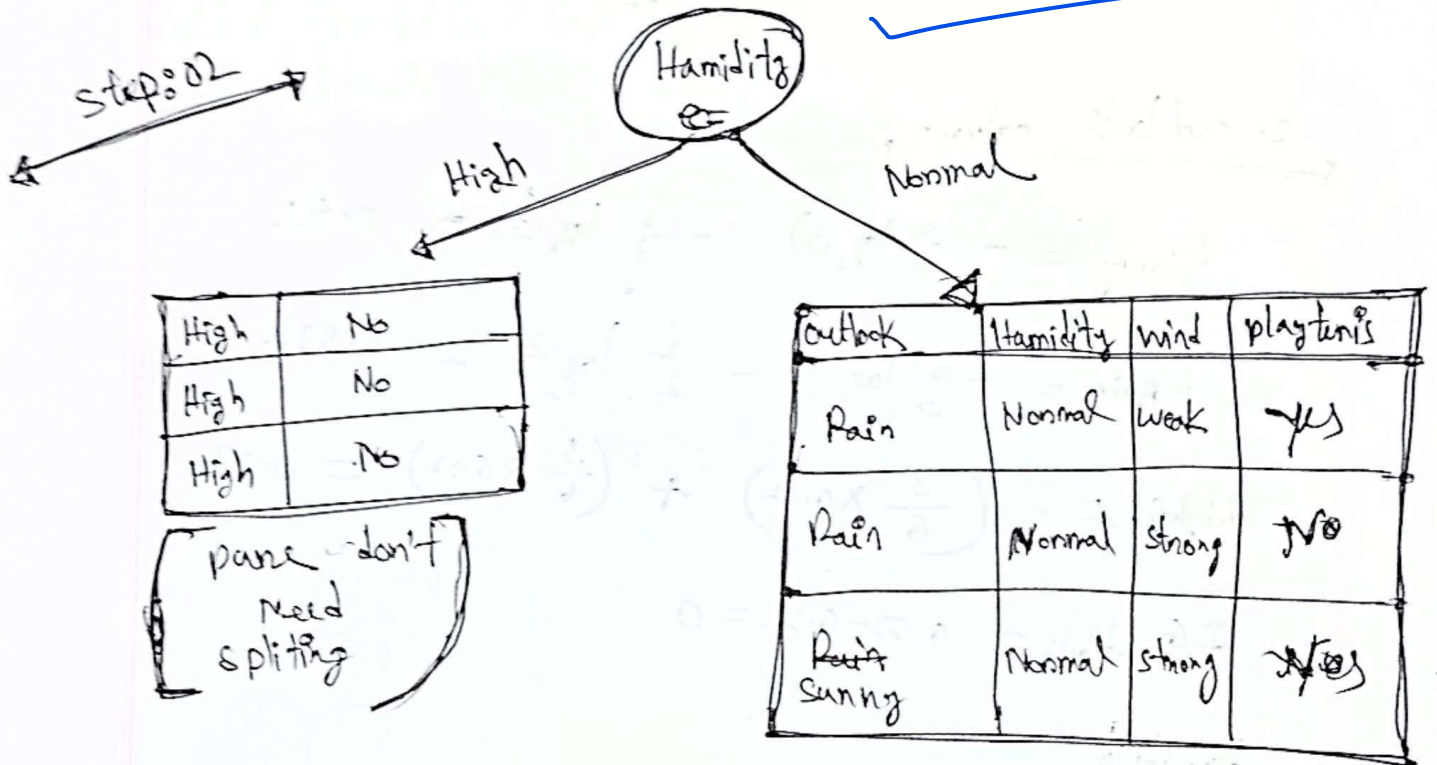$$E_{weak} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$E_{No} = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) = 0.811$$

$$\text{Weighted } E = \frac{2}{6} \times 1 + \frac{4}{6} \times 0.811 = 0.874$$

$$IG_{wind} = 0.92 - 0.874 = 0.046$$

CamScanner

$$IG_{Humidity} > IG_{wind} > IG_{outlook}$$

∴ Root Node (Humidity)

Step:02

Hamidity

High

Nonimal

| High | No |
|------|-----|
| High | No |
| High | No |

( pane don't
Need
spliting )

| outlook | Hamidity | wind | play tenis |
|---------|----------|------|------------|
| Rain | Nonmal | weak | yes |
| Rain | Nonmal | strong | No |
| Rain Sunny | Nonmal | strong | Yes |

$$E_{Nonimal} = 0.92 \quad (\text{from previous})$$

**Wind:**

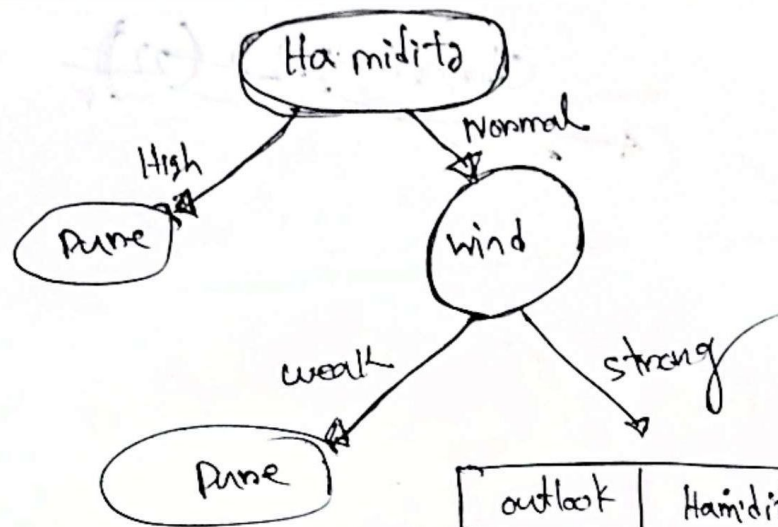$$E_{strong} = 0 - \frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$E_{weak} = 0$$

$$\text{weighted entropy} = \left(\frac{1}{3} \times 0\right) + \left(\frac{2}{3} \times 0\right) = 0.33$$

$$IG_{wind} = 0.92 - 0.33 = 0.92 \quad 0.59$$

( highest IG Grain) [ split beme on this. ]

values
may be
same
take
anyone

**Diagram 1 (top):**

Humidity
- High → Pure
- Normal → Wind
  - weak → Pure
  - strong →

| outlook | Hamidity | wind | play |
|---------|----------|-------|------|
| Rain | Normal | strog | No |
| Sunny | Normal | Normal | Yes |

**Diagram 2 (bottom):**

Hamidity
- High → Pure
- Normal → Wind
  - weak → Pure
  - strong → outlook
    - Pure
    - Pure