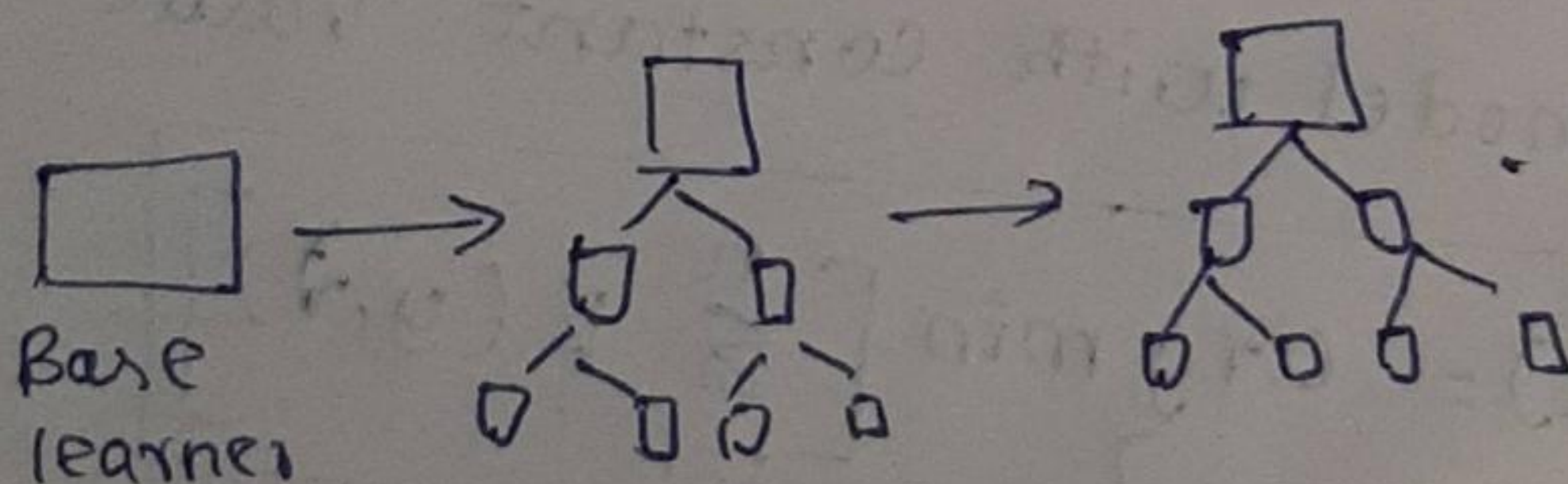# Gradient Boosting:-

* It is the numerical optimisation problem where the objective is to minimise the loss function of the model by adding weak learners using gradient descent.

* Gradient descent is a first order iterative optimisation algorithm for ~~finding~~ finding local minimum of differentiable function.

* Gradient boosting is based on minimising a loss function.

* The basic idea is it will first set a basic learner and then the decision trees will get added as Per the ~~hy~~ value of hyperparameter.

* This is worked out by taking the previous model's residuals as the i/p to the present model.



Base learner

* The main difference with adaboost is adaboost focuses on misclassified observation where as gradient boost trains learners based upon minimising the loss function of learners that is training ~~as~~ on the residuals of the model.

Dataset :- ~~I/P~~ Indep | O/P

| Exp | degree | Salary (in K) |
|-----|--------|---------------|
| 2 | BE | 50K |
| 3 | PHD | 70 K |
| 4 | BE | 60K |

Input to gradient boosting:-

① $\{x_i, y_i\}$

$x_i \Rightarrow$ independent features

$y_i \Rightarrow$ dependent features

② $d(y, F(x))$

↳ Loss function like mse, RMSE.

(Loss function should be differentiable)

③ No. of trees needed.

Pseudo Algorithm:-

① Initialize model with constant value

$$F_0(x) = \arg\min_\gamma \left[ \sum_{i=1}^n L(y, \gamma) \right]$$

where $L(y, \gamma)$ is Loss function

$y$ is data pts

$\gamma$ is predicted values

Now, defining a loss function

$$Loss = \sum_{i=1}^{n} \frac{1}{2}(y - \hat{y})^2$$
$$L(y, \hat{y})$$

* we want to find $\hat{y}$ in such a way that the loss function should be reduced.

& Substituting values from dataset in loss function

$$= \frac{1}{2}(50 - \hat{y})^2 + \frac{1}{2}(70 - \hat{y})^2 + \frac{1}{2}(60 - \hat{y})^2$$

, v To find the $\hat{y}$ in order to minimise loss function we want to find first order derivative. ~~this step~~ actually this step uses gradient descent.

$$\overset{0}{\underset{\substack{\text{critical} \\ \text{Pt. so it is 0}}}{\downarrow}} = \frac{\partial}{\partial \hat{y}}\left[\frac{1}{2}(50 - \hat{y})^2 + \frac{1}{2}(70 - \hat{y})^2 + \frac{1}{2}(60 - \hat{y})^2\right]$$

$$0 = \frac{2}{2}(50 - \hat{y})(-1) + \frac{2}{2}(70 - \hat{y})(-1) + \frac{2}{2}(60 - \hat{y})(-1)$$

$$0 = -50 + \hat{y} - 70 + \hat{y} - 60 + \hat{y}$$

$$0 = -180 + 3\hat{y}$$

$$3\hat{y} = 180$$

$$\boxed{\hat{y} = 60}$$

$\downarrow$

$$\boxed{\text{Constant value} = 60}$$

* This is $\hat{y}$ for the Base learner.

updating dataset with $\hat{y}$

| exp | degree | salary (in k) | $\hat{y}$ |
|-----|--------|---------------|-----------|
| 2   | BE     | 50k           | 60        |
| 3   | PHD    | 70k           | 60        |
| 4   | BE     | 60k           | 60        |

② Iterate $m = 1$ to $m$ (no. of trees)

i) compute pseudo residuals (Pseudo error)

$$r_{im} = -\left[\frac{\partial h\,(y,\,F(x_i))}{\partial F(x_i)}\right].$$

This is nothing but

$$Loss = \frac{1}{2}(y-\hat{y})^2$$

$$= \frac{\partial}{\partial \hat{y}}\left(\frac{1}{2}(y-\hat{y})^2\right)$$

$$\frac{\partial h}{\partial \hat{y}} = \frac{2}{2}(y-\hat{y})(-1)$$

$$\frac{\partial h}{\partial \hat{y}} = -(y-\hat{y})$$

$$\boxed{-\frac{\partial h}{\partial \hat{y}} = y-\hat{y}}$$

So,

$$r_{im} = -\left[\frac{\partial h\,(y,\,F(x_i))}{\partial F(x_i)}\right] = -\frac{\partial h}{\partial \hat{y}} = y-\hat{y}$$

* Basically $r_{im}$ is o/p feature — residual $(\hat{y})$

of Base learner. ($r_{im}$ means residual of Base model)

updating dataset with $r_{im}$

| exp | degree | salary (in k) | $\hat{y}$ | $r_{im}$ $(y-\hat{y})$ |
|-----|--------|---------------|-----------|------------------------|
| 2   | BE     | 50K           | 60        | -10                    |
| 3   | PHD    | 70K           | 60        | 10                     |
| 4   | BE     | 60K           | 60        | 0                      |

ii) Fit a base learner $h_m(x)$ with $\{x_i, r_{im}\}$

where $x_i$ = Independent features

(Training decision tree)

$r_{im}$ (residuals) = dependent feature

* This is a decision tree regressor

iii) Finding $\mathcal{V}_m$ in order to minimise the loss function.

$$\mathcal{V}_m = \arg\min_{\mathcal{V}} \sum_{i=1}^{n} L\left(y_i, F_{m-1}(x_i) + \mathcal{V}\right)$$

This term clearly states its using output of previous iteration
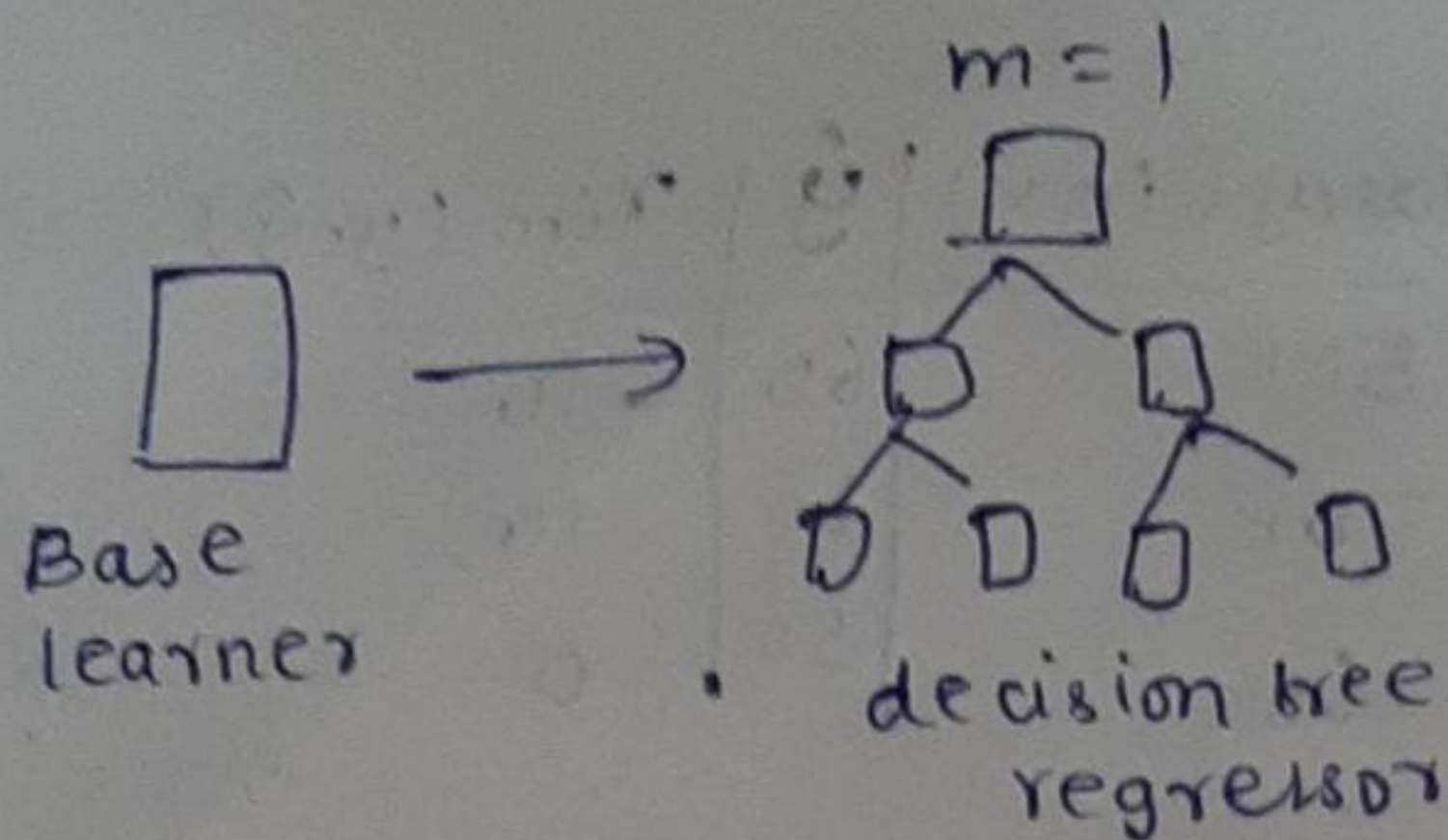
* This $\mathcal{V}_m$ same as $F_0(x)$.

Previous model value added.

$$\mathcal{V}_m = \sum_{i=1}^{n} \frac{1}{2} \left(y_i - (60 + \hat{y})\right)^2$$

IV) update model

$$F_m(x) = F_{m-1}(x) + \alpha(h(x))$$

m=1



Base learner → decision tree regressor

$\alpha \Rightarrow$ learning rate

usually b/w

0 to 1

$$F_1(x) = F_{(1-1)}(x) + \alpha(h(x))$$

$\rightarrow \gamma_{11}$

$$= F_0(x) + \alpha(h(x))$$

(with only one record)

First time

$$= 60 + 0.1(-10)$$

$$= 60 - 1.0$$

$y_1 = 50k$ & huge diff

$F_1(x) = 59$ So, goes to next tree

$$= 59.0$$

* It will calculate for all the records.

& and if the diff is huge, it will go to next tree.

* usually decision trees leads to overfitting. So using less no. of trees will be ideal or else we can use regularization techniques to reduce generalize.