

ENSEMBLED METHODS

ADA BOOST

Mathematical Intuitions :-

- ✓ Boosting is an ensembled Technique where we create sequentially.
- ✓ combined multiple models
- ✓
 - ① Ada Boost
 - ② Gradient Boost
 - ③ XGB

Ada Boost:-

Bagging \Rightarrow parallel creation of models
 ✓ Models are independent

Bootstrap \rightarrow Aggregation

Eg:- Random Forest

✓ Sequential implementation \Rightarrow $y_{\text{Pred}} = DT_1\alpha_1 + DT_2\alpha_2 + DT_3\alpha_3 + \dots + DT_n\alpha_n$

✓ Additive modelling \rightarrow addition of all models

✓ models are dependent on each other

$$\text{Ada Boost} := h(\theta) = h_1(\theta)\alpha_1 + h_2(\theta)\alpha_2 + h_3(\theta)\alpha_3 + \dots + h_n(\theta)\alpha_n$$

- ✓ Classification
- ✓ Regression

$h(\theta)$ = function \leftarrow Decision Tree

α = weight



I Stump \rightarrow Weak learner

Training Accuracy \rightarrow underfitting

Wrong classification

$w_1 + w_2 + w_3 \dots = \text{Strong Learner}$
 Low bias (Overfitting)
 High variance
 High bias (Underfitting)
 Low variance

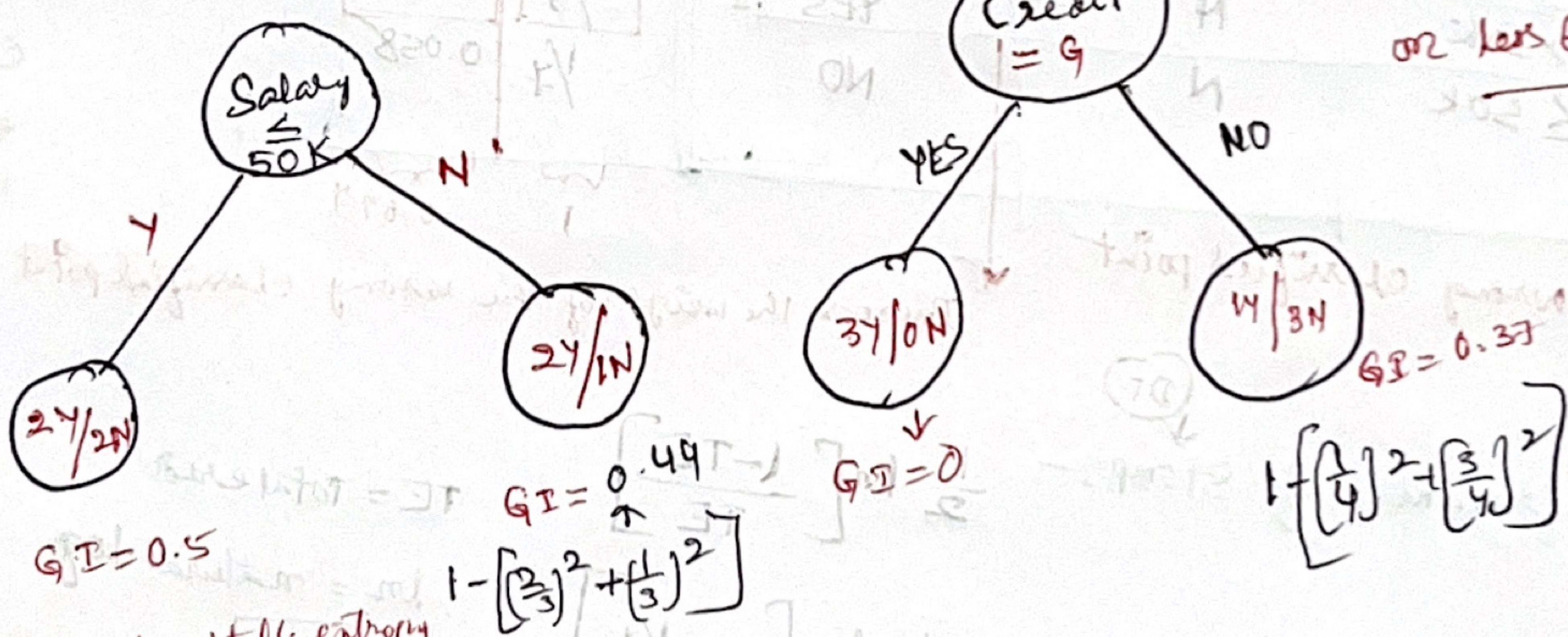
bias variance $\left\{ \begin{array}{l} \text{Boosting} \\ \text{Boosting} \end{array} \right.$

Data Set

<u>Salary</u>	<u>credit Score</u>	<u>Approval</u>
$\leq 50K$	B	NO
$\leq 50K$	G	YES
$\leq 50K$	G	YES
$> 50K$	B	NO
$> 50K$	N	YES
$> 50K$	G	YES
$\leq 50K$	N	NO

Step 1:

Create a Decision Tree



Step 2: Find out suitable entropy

(Algorithms can be used)

$ID_3 = \text{Entropy} \rightarrow \text{Information Gain}$

CART = Gini Impurity

Gini Impurity \Rightarrow

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

$$\left[\frac{4}{9} + \frac{1}{6} \right] = \frac{5}{9}$$

$$\frac{4+1}{9} = \frac{5}{9}$$

$$\left[\frac{1}{16} + \frac{9}{16} = \frac{10}{16} \right]$$

$$1 - \frac{10}{16} = \frac{6}{16} = \frac{3}{8}$$

$$1 - \frac{5}{9} = \frac{9-5}{9} = \frac{4}{9}$$

Step 3:-

① Decision Tree with maximum one Layer is called 'Stump' It is also known as Depth 1.

② Weight - Assign Sample weight

Salary	Credit Score	Target column Approval	n = no of rows Sample Weight ($\frac{1}{n}$)	
			updated weight	
<=50K	B	NO	$\frac{1}{7}$	0.058
<=50K	G	YES	$\frac{1}{7}$	0.058
<=50K	G.	YES	$\frac{1}{7}$	0.058
> 50K	B	NO	$\frac{1}{7}$	0.058
> 50K	G	YES	$\frac{1}{7}$	0.058
> 50K	N	YES	$\frac{1}{7}$	0.349
≤ 50K	N	NO	$\frac{1}{7}$	0.058

→ wrong classified point

Performance of Stump:-

$$\frac{1}{2} \ln \left[\frac{1 - TE}{TE} \right]$$

TE = total error

$$= \frac{1}{2} \ln \left[\frac{1 - \frac{1}{7}}{\frac{1}{7}} \right]$$

ln = natural log

Log_e = log with base e

$$\text{performance of Stump} = \frac{1}{2} \ln \left[\frac{6}{7} \right] \approx 0.896$$

$$\alpha_1 = 0.896$$

→ Now increase the weight of wrong classified point and decrease the weight of correctly classified points.

② Calculate the updated weight

$$\text{Weight} \times e^{-\text{Performance of Stump}}$$

formula for
correct classified
points

$$\Rightarrow \frac{1}{7} \times e^{-(0.896)}$$

$$\Rightarrow 0.058$$

Performance of Stump

for wrongly
classified
points

$$= \frac{1}{2} + e^{(0.896)}$$

$$\approx 0.349$$

≈ 0.349 2011-08-03 Function assign some basis
Each record assign some basis

Frost, Lida S. P. Oldsmar, Fla.

Normalized

weights

0.083

0.083

0.083

0.08
22

0.50

0.083

100

$$\text{LPI} + \dots + \text{GPI} \dots 0.16 - 0.94$$

$$0.24 - 0.32$$

$$0.32 - 0.40$$

$$\boxed{0.40 - 0.90}$$

$$0.90 - 0.96$$

- first bin size
- second bin size

Cumulative Sum

Based on bucket size only we will select wrong classified points
Range = $[0-1] \Rightarrow \{0.2, 0.4, 0.5, 0.6, 0.7\}$

How to Select points for 2nd DT

Data set 2

Data Set 2		
Salary	Credit Score	Approve
<=50K	G	N
> 50K	N	Y
> 50K	N	Y
> 50K	N	Y

101
Here chances for selection of wrong
classified points is higher

$$\frac{0.896}{\alpha_1} + \frac{0.650}{\alpha_2} + \frac{0.24}{\alpha_3} + \frac{0.30}{\alpha_4} (\text{NO})$$

↑
10

$$1136 (\text{YES}) + 0.350 (\text{NO})$$

Performance of Stump? :-

$$1.136 \text{ (yes)} + 0.350 \text{ (No)}$$

$$0|8 = yes$$

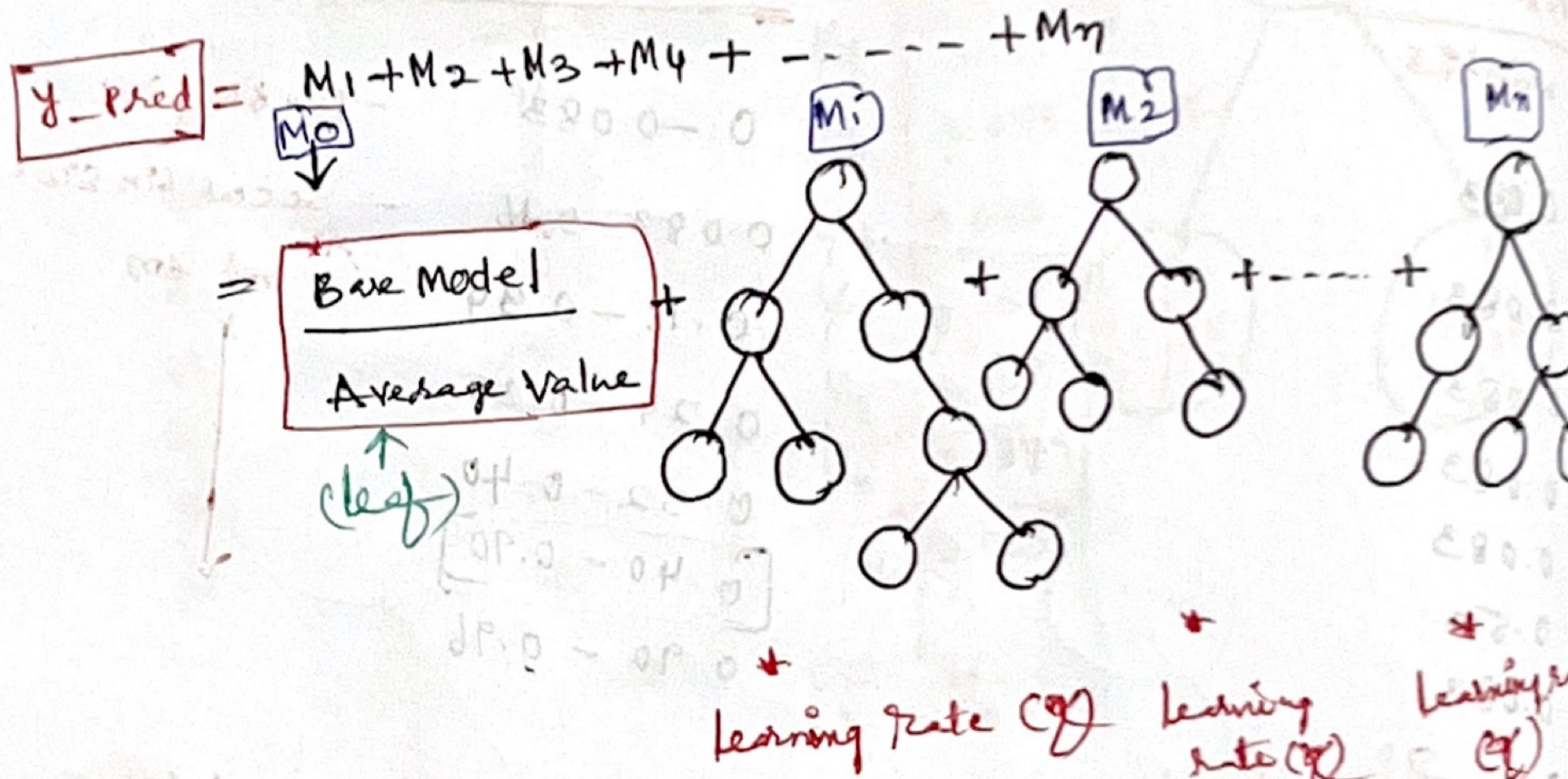
↑ additive learning

How many Decision Trees \rightarrow Hyperparameter $\rightarrow n_{estimators} = 50$ (Default)
 $estimator = \text{model}$
Learning rate

GRADIENT BOOSTING

(Classification & Regression)

Sequence of the Model



Data Set

Iqr CGPA

Salary Admissions

Prediction

		82	66
90	8	62	6.6
100	9	52	6.6
110	6	102	6.6
120	9	42	6.6
80	5	8 + 6 + 5 + 10 + 9 5	6.6

$$\begin{aligned} & (\mathbf{x}_i, y_i)_{i=1}^n \\ & \text{IP features} \\ & L(y, p(x)) \\ & \text{Loss Function} \end{aligned}$$

1. Initialize model with a constant value
 $f_{0,0} = \arg \min_r \sum_{i=1}^n L(y_i, r)$

$$\text{Residual} = \text{Loss}(\text{Red}) = (y - \text{pred})$$

$$= 8 - 6.6 = 1.4$$

$$6 - 6.6 = -0.6$$

$$5 - 6.6 = -1.6$$

$$10 - 6.6 = 3.4$$

$$4 - 6.6 = -2.6$$

1.4
-0.6
-1.6
3.4
-2.6

Gradient Boosting Regression

In [24]:

```
import pandas as pd
df=pd.DataFrame({
    'iq':[90,100,110,120,80],
    'cgpa':[8,7,6,9,5],
    'salary':[3,4,8,6,3]
})
df.head()
```

Out[24]:

	iq	cgpa	salary
0	90	8	3
1	100	7	4
2	110	6	8
3	120	9	6
4	80	5	3

In [25]:

```
import numpy as np
base_model=np.mean(df['salary'])
```

In [26]:

```
df['residual1']=df['salary']-base_model
```

In [27]:

```
df['residual1']
```

```
Out[27]: 0    -1.8  
1    -0.8  
2     3.2  
3     1.2  
4    -1.8  
Name: residual1, dtype: float64
```

```
In [28]: from sklearn.tree import DecisionTreeRegressor  
model=DecisionTreeRegressor()
```

```
In [29]: df['base_model']=base_model
```

```
In [30]: df
```

```
Out[30]:   iq  cgpa  salary  residual1  base_model  
0    90     8      3      -1.8        4.8  
1   100     7      4      -0.8        4.8  
2   110     6      8      3.2         4.8  
3   120     9      6      1.2         4.8  
4    80     5      3      -1.8        4.8
```

```
In [31]: from sklearn.model_selection import train_test_split  
X_train,X_test,y_train,y_test=train_test_split(df[['iq','cgpa']],df.residual1,test_size=0.2,random_state=42)
```

```
In [32]: model.fit(X_train,y_train)
```

```
Out[32]: ▾ DecisionTreeRegressor  
DecisionTreeRegressor()
```

```
In [33]: model.score(X_train,y_train)
```

```
Out[33]: 1.0
```

```
In [34]: from sklearn import tree
import matplotlib.pyplot as plt
#fig,ax=plt.subplots(figsize=(10,5))
plt.figure(figsize=(10,5))
tree.plot_tree(model, filled=True)
plt.show()
```

