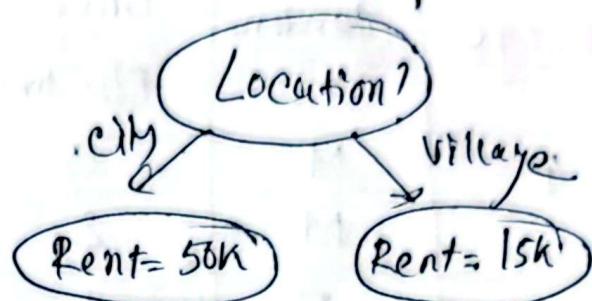


Regression Tree



When a decision tree classifies things into categories — **Classification Tree**

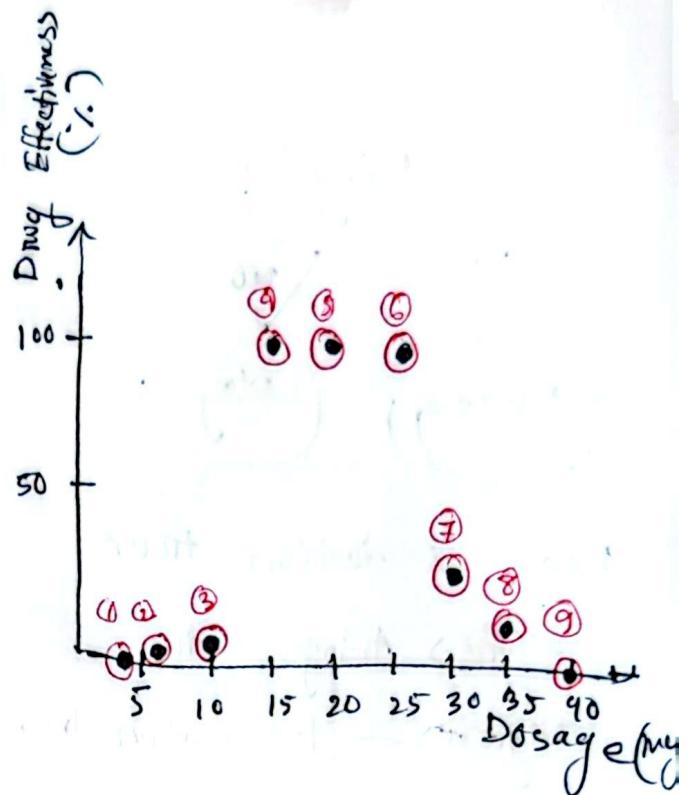
When decision tree predicts numeric values — **Regression tree**.

If we have a dataset with Drug "dosage" feature and Drug "Effectiveness" as target, which plots like the right picture, then we can easily fit a best fit line. Cause here ~~tree~~ when dosage increases, the drug effectiveness also increases. Then for any drug dosage = x , we can easily predict the effectiveness of the drug = y from this line. This is what we do in Linear Regression.

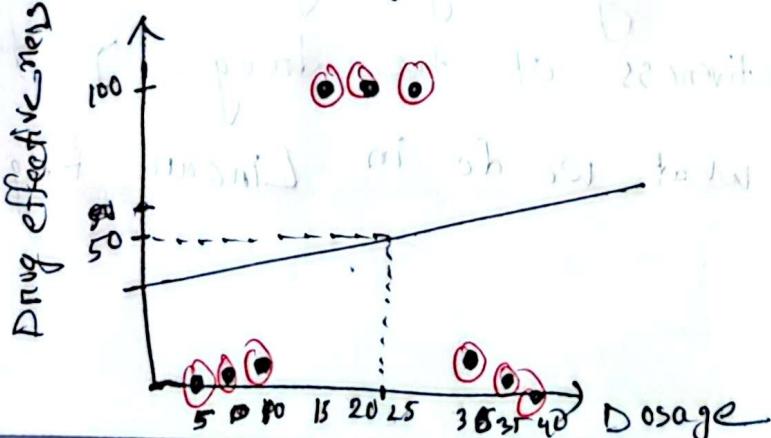
Dosage (X)	Drug Effectiveness (Y)
1	1.5
2	2.5
3	3.5
4	4.5
5	5.5
6	6.5
7	7.5
8	8.5
9	9.5
10	10.0

But if we have a dataset like below →

	Dosage	Gender	Drug Effectiveness
①	9	M	0
②	6	M	2
③	10	F	3
④	15	F	100
⑤	20	M	100
⑥	25	M	100
⑦	30	F	20
⑧	35	F	10
⑨	40	M	0

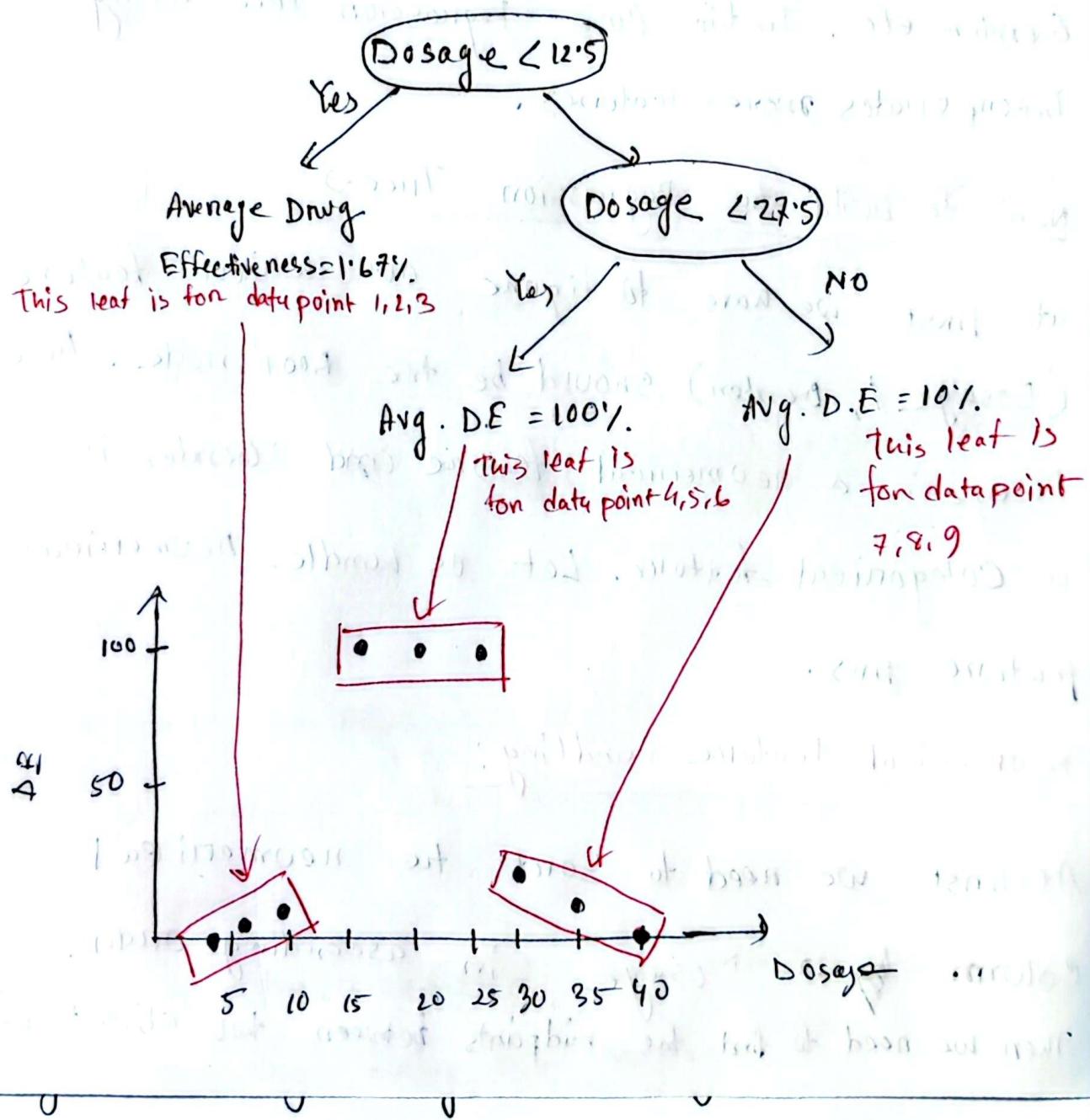


The plot does not give us an uniform cluster where we can say when the drug dosage increases the drug effectiveness also increases. Here low dosage → not effective, Moderate dosage → works well, and high dosage → works moderate to does not work at all. If we fit a straight line here →



The straight line will not be able to predict drug effectiveness. Hence if we predict for dosage = 23 mg, we get effectiveness = 50% (approx). But we know moderate dosage should give us 100% effectiveness.

So, we can use Regression Tree here to make prediction
 as the target column is numerical. Each leaf of the regression tree predicts numeric values.



Now if we want to predict for dosage = 23 mg, from the decision tree we get the drug effectiveness = 100%.

Even though we can predict this values only by looking the graph/plot, but it gets much more complicated when there are more features in the dataset like Gender etc. In this case Regression tree easily incorporates more features.

Now to build the Regression Tree

at first we have to figure out which feature (Dosage & Gender) should be the Root node. Here Dosage is a numerical feature and Gender is a categorical feature. Let us handle numerical feature first.

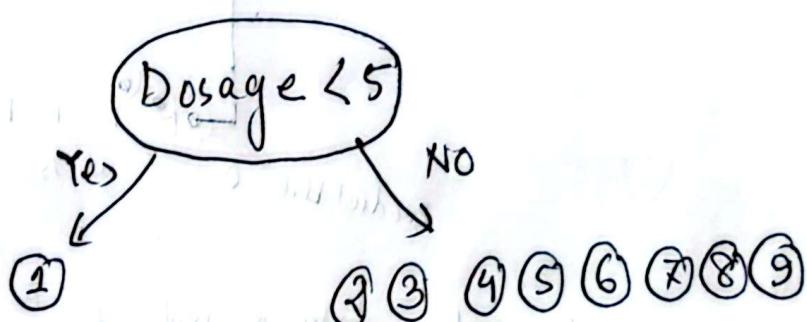
Numerical Feature handling:

At first we need to sort the numerical column \rightarrow Dosage, in ascending order. Then we need to find the midpoints between two adjacent points.

$$\begin{array}{ccccccccc} \text{Dosage} \rightarrow & 4 & 6 & 10 & 15 & 20 & 25 & 30 & 35 & 40 \\ & \downarrow \\ \text{Mid points} \rightarrow & 5 & 8 & 12.5 & 17.5 & 22.5 & 27.5 & 32.5 & 37.5 \end{array}$$

Now we have to set every mid point as a threshold for the dosage column to see how well it splits the data. To measure which threshold splits better, we will calculate $SSR = \text{sum of square residual}$.

Let us take the first mid point 5 as threshold for the root node.



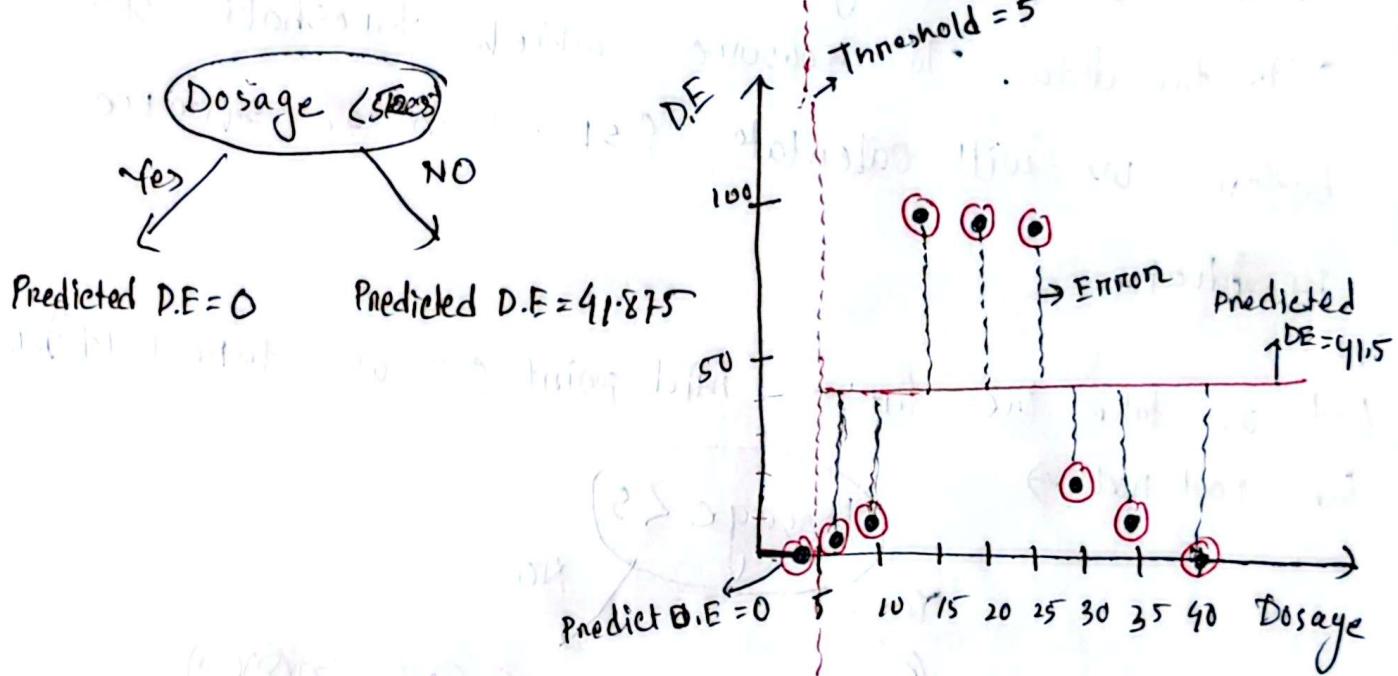
only the first data point is less than 5, so it comes to the left leaf. all other data points goes to the right branch.

We calculate the **Predicted** drug effectiveness for every branch by taking the average of the

drug effectiveness of data points in each branch.

left branch \rightarrow data point (1) \rightarrow average $D.E = \frac{0}{1} = 0$

Right " \rightarrow data point (2, 3, 4, 5, 6, 7, 8, 9) \rightarrow average $D.E = \frac{2+3+100+100+100+20+10+0}{8} = 41.875$



In the left branch only data point 1 is present, which has $actual D.E = 0$. The predicted $D.E = 0$ as well.

If we calculate error between Actual D.E and predicted D.E using SSR we get \rightarrow

$$SSR_{left\ leaf} = (0 - 0)^2 = 0$$

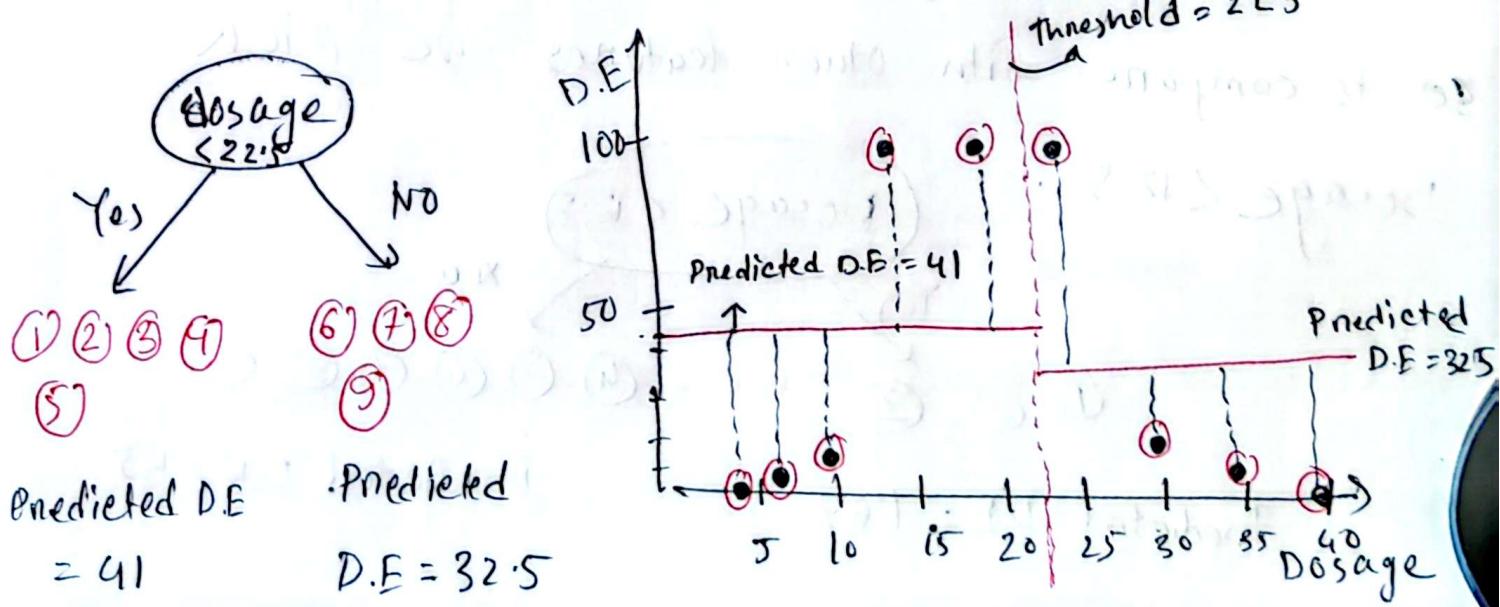
Actual D.E Predicted D.E

In the right branch we have all other data points, and the actual drug effectiveness are $\Rightarrow \textcircled{1} \rightarrow 2$, $\textcircled{2} \rightarrow 3$, $\textcircled{4} \rightarrow 100$, $\textcircled{5} \rightarrow 100$, $\textcircled{6} \rightarrow 100$, $\textcircled{7} \rightarrow 20$, $\textcircled{8} \rightarrow 10$, $\textcircled{9} \rightarrow 0$. But the predicted D.E = 41.875. There are huge errors present between actual D.E and predicted D.E, which we can calculate using SSR.

$$\begin{aligned}
 \text{SSR}_{\text{Right}} &= \text{Actual} - \text{Predicted} \\
 &= (2 - 41.875)^2 + (3 - 41.875)^2 + (100 - 41.875)^2 + \\
 &\quad (100 - 41.875)^2 + (100 - 41.875)^2 + (20 - 41.875)^2 \\
 &\quad + (20 - 41.875)^2 + (0 - 41.875)^2 \\
 &= 16485
 \end{aligned}$$

$$\text{So, } SSR(\text{Dosage} < 5) = SSR_{\text{left}} + SSR_{\text{right}} \\ = 0 + 16485 \\ = 16485.$$

Now, let us calculate SSR for dosage ≤ 22.5



$$SSR(\text{Dosage} \leq 22.5) = \left[(0 - 41)^2 + (2 - 41)^2 + (3 - 41)^2 + (100 - 41)^2 \right]$$

left branch

$$+ (100 - 41)^2$$

$$\left. \begin{array}{l} \\ \\ \\ \end{array} \right\} + (100 - 32.5)^2 + (120 - 32.5)^2 + (10 - 32.5)^2$$

right branch

$$+ (0 - 32.5)^2$$

$$= 17883$$

If we calculate SSR for all other threshold we get,

$$\text{Dosage} \leq 5$$

$$16485$$

$$\text{Dosage} \leq 22.5$$

$$\text{Dosage} \leq 8$$

$$14670$$

$$\text{Dosage} \leq 27.5$$

$$\text{Dosage} \leq 12.5$$

$$12355$$

$$\text{Dosage} \leq 32.5$$

$$17172$$

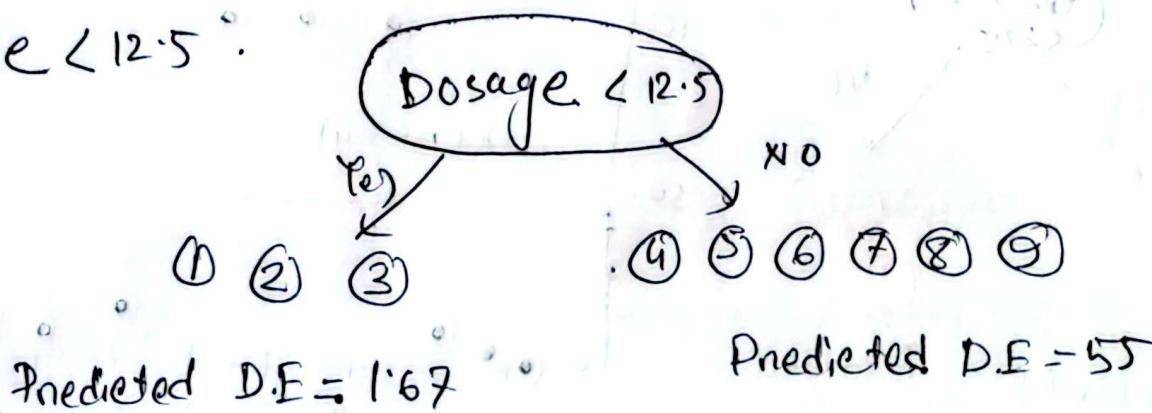
$$\text{Dosage} \leq 37.5$$

$$17883$$

$$14709$$

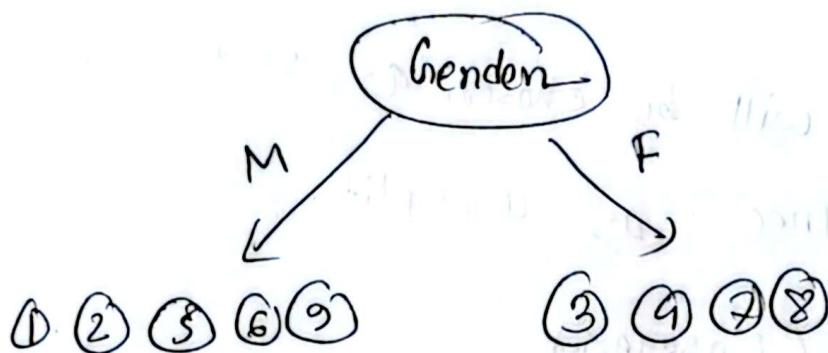
$$15374$$

$$16485$$



Categorical Feature handling:

handling categorical feature is easy as we can just split the datapoint according to the categories of this feature like below.

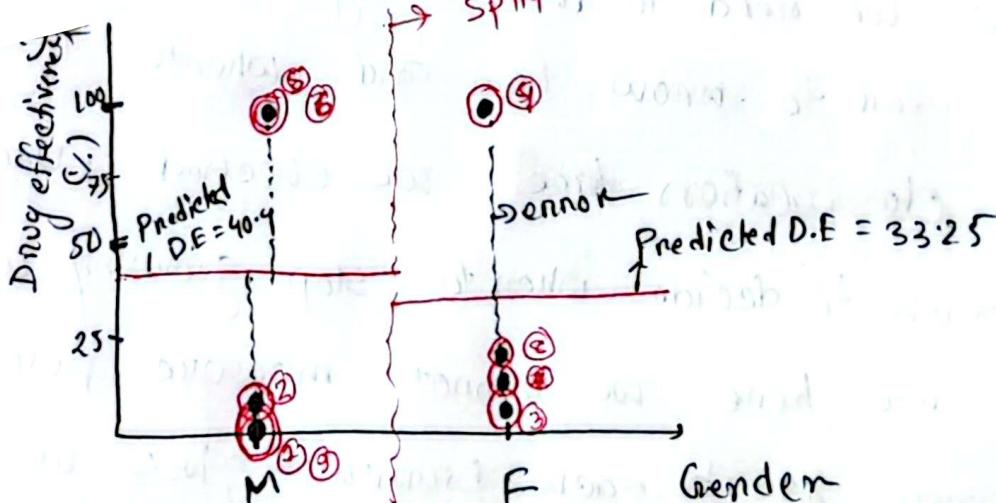


Average of Drug effectiveness
of these points →

$$\text{Predicted D.E.} = \frac{0+2+100+100+0}{5} = 40.4$$

Average of Drug effectiveness
of these data points →

$$\text{Predicted D.E.} = \frac{3+100+20+10}{4} = 33.25$$



$$\therefore \text{SSR}(\text{Gender}) = (0 - 40.4)^2 + (2 - 40.4)^2 + (100 - 40.4)^2 + (100 - 40.4)^2 + (0 - 40.4)^2$$

$$\text{Gender} = M$$

$$+ (3 - 33.25)^2 + (100 - 33.25)^2 + (20 - 33.25)^2 + (10 - 33.25)^2$$

$$\text{Gender} = F$$

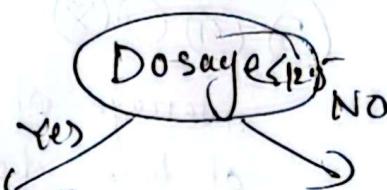
$$SSR(\text{Gender}) = 17929.25$$

SSR of all features :

$$SSR(\text{Dosage} < 12.5) = 12.355$$

$$SSR(\text{Dose} \neq \text{Gender}) = 17929.25$$

∴ Dosage < 12.5 will be chosen as root node of the regression Tree as a split.



$$\text{Predicted D.E} = 1.67$$

$$\text{Predicted D.E} = 55$$

① ② ③

④ ⑤ ⑥ ⑦ ⑧ ⑨

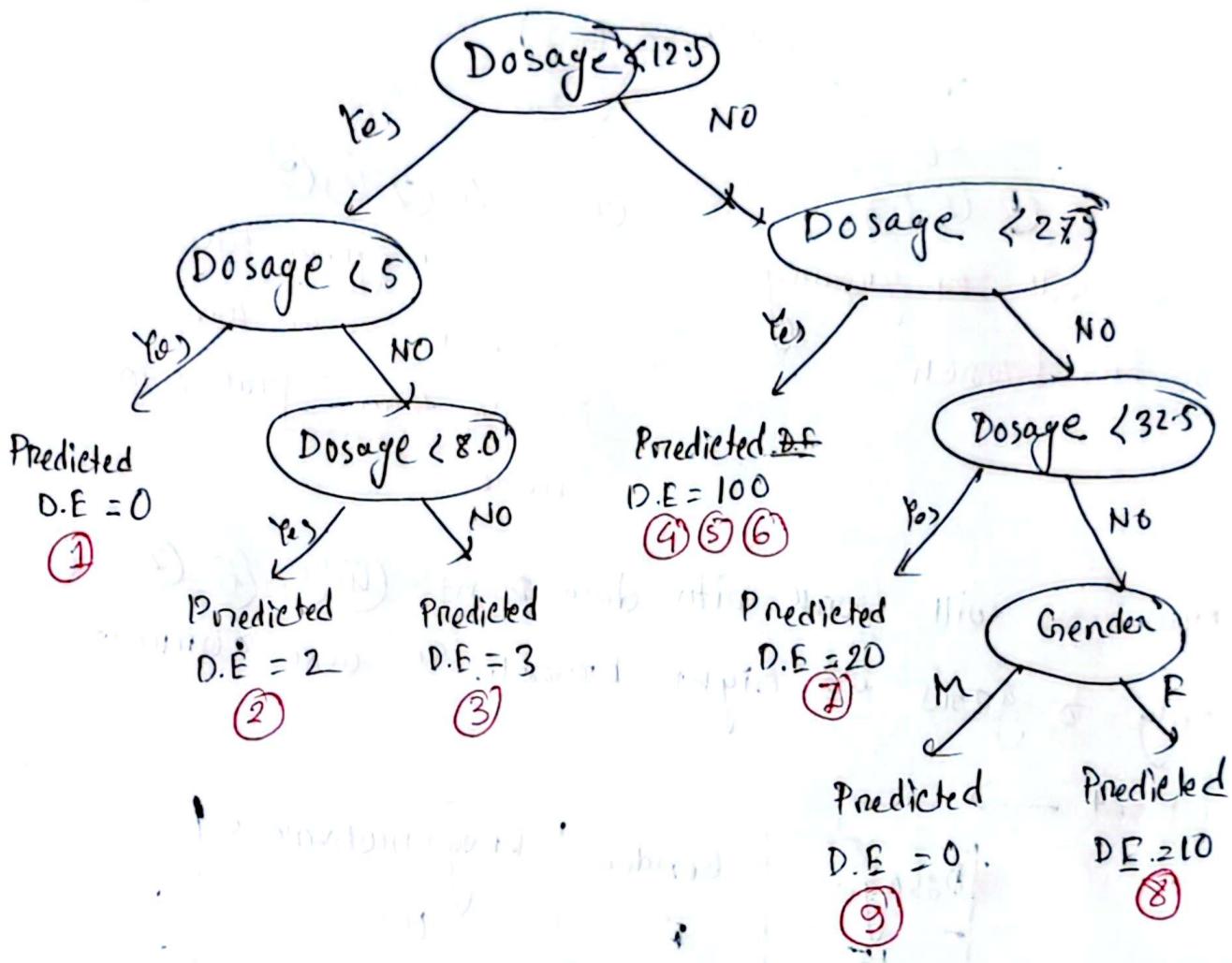
Now we need to further grow this tree. But we need to know how and when to stop.

In classification tree we checked purity of a branch to decide when to stop growing a tree.

However here, we cannot measure purity of a branch as each branch gives us a numerical prediction, rather than any class.

So to know this. Moreover if we keep growing

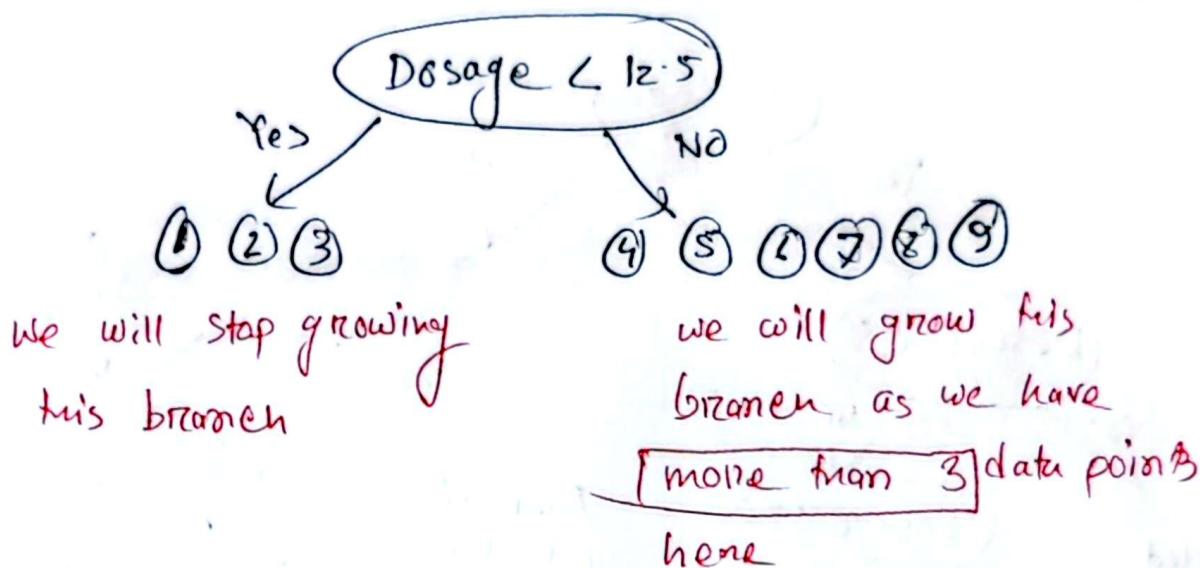
We might end up with a tree like this →



Hence we can see, each leaf ends up with one sample / data point (except 'Yes' branch of 27.5 as $④ ⑤ ⑥$ has $D.E = 100$). So, this regression tree has overfitted into our training dataset and it will not work well for unknown data.

So To solve this issue, like the classification tree we can set limit for leaf nodes size, "Do not split further if size is 3 or less than 3 data point in a leaf". If we enforce this limit

in this tree the ^{2nd} level of the tree →



Now we will work with data point ④⑤⑥⑦⑧⑨ only to grow the right branch. So our current dataset →

Dosage	Gender	Drug effectiveness
15	F	100
20	M	100
25	M	100
30	F	20
35	F	10
40	M	0

Again we have to calculate $SSR(\text{Gender})$ and $SSR(\text{Dosage} < \text{threshold})$ and compare those to decide the node of this branch. As Gender was not used as node before and dosage feature can have different thresholds, we will consider both.

Calculating SSR for Dosage feature

Let us sort the dosage feature again and find the midpoints →

Dosage →	100	20	25	30	35	40
	↓	↓	↓	↓	↓	↓
	27.5	22.5	27.5	32.5	37.5	

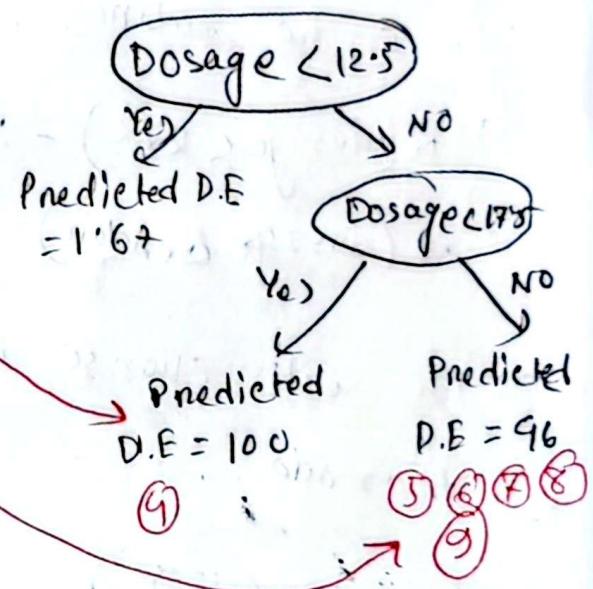
Now we have to calculate SSR by setting all of these midpoints as threshold.

Let us calculate for →

SSR (Dosage < 17.5)

$$= \frac{1}{6} [(100 - 100)^2 + (100 - 46)^2 + (100 - 46)^2 + (20 - 46)^2 + (30 - 46)^2 + (0 - 46)^2]$$

$$= 9920$$



Similarly if we calculate for all of the midpoints we get →

Dosage < 17.5
↓
9920

Dosage < 22.5
↓
6275

Dosage < 27.5
↓
200 → lowest

Dosage < 32.5
↓
7425

Dosage < 37.5
↓
13076

calculating SSR for Grenden feature \rightarrow

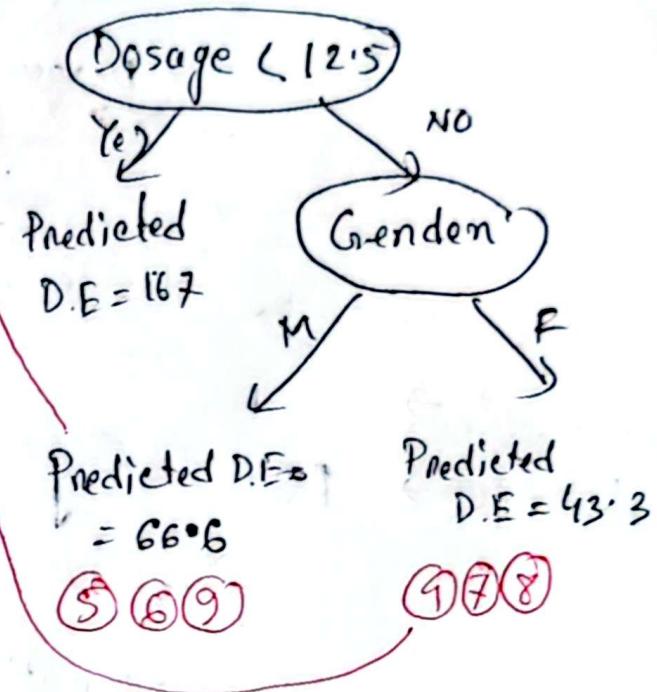
$$SSR(\text{Grenden}) =$$

$$= (100 - 66.6)^2 + (100 - 66.6)^2 + (8 - 66.6)^2$$

$$+ (100 - 43.3)^2 + (20 - 43.3)^2 + (10 - 43.3)^2$$

$$= 6666.68 + 4866.67$$

$$\therefore 11533.35$$



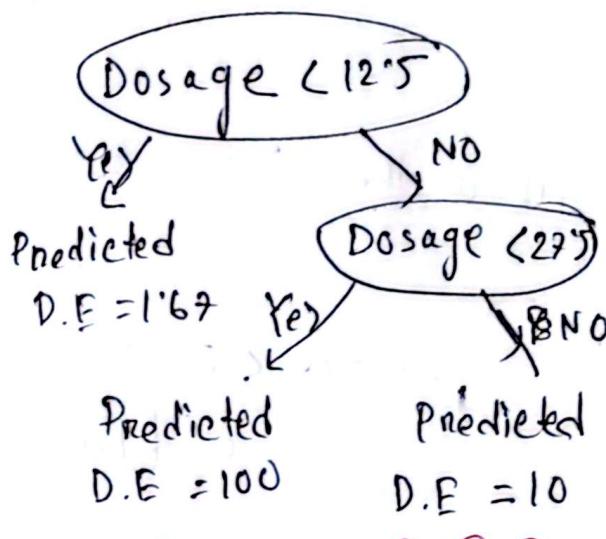
\therefore SSR of the features for third level node \rightarrow

$$SSR(\text{dosage} < 27.5) = 200 \rightarrow \text{Lowest}$$

$$SSR(\text{dosage Grenden}) = 11533.35.$$

so, we will choose lowest SSR, dosage < 27.5.

as thenode..



④ ⑤ ⑥ ⑦ ⑧ ⑩
3 datapoint 3 datapoint

As the leaf nodes have 3 data points each, which is our limit, we will not grow the tree further.