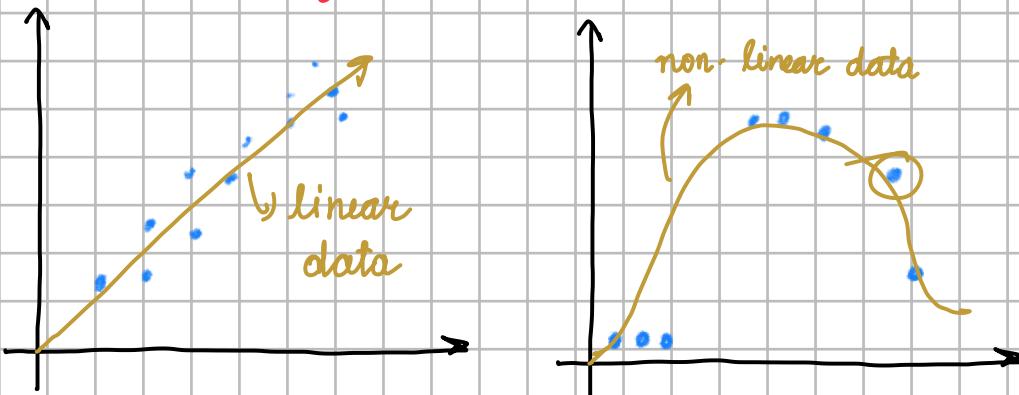


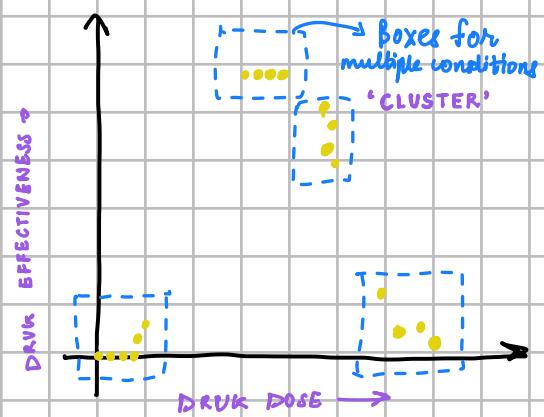
# REGRESSION TREE (FROM SCRATCH)

Can handle non-linear data much better than the other regression methods.



Regression tree is used here.

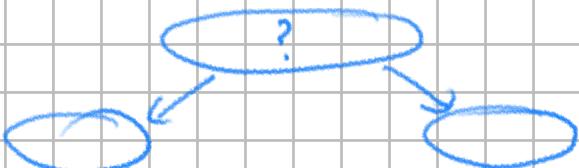
## \* Visualization



Why does this tree makes good prediction??

Problem: We have a training dataset which consists of drug effectiveness of different dosage. Given this training data, we want to build a regression tree that uses drug doses to predict drug effectiveness.

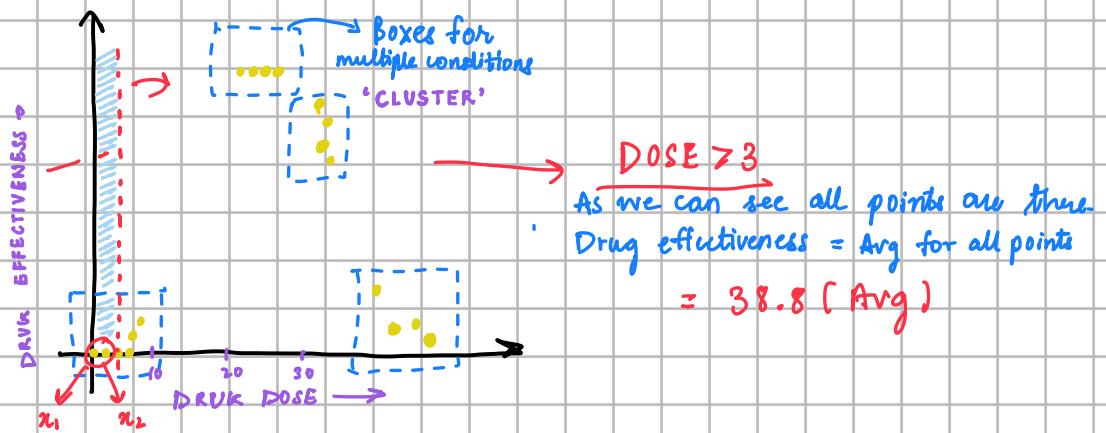
Step ①: We need to decide root node first!



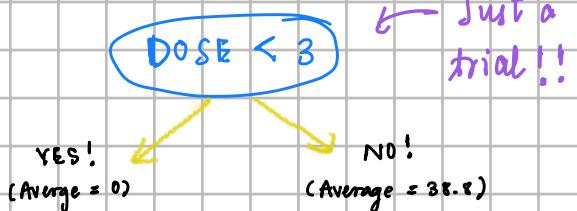
Step ② calculate the avg. of first two doses (equals to 3)

$$\frac{x_1 + x_2}{2} = 3 \text{ (avg)}$$

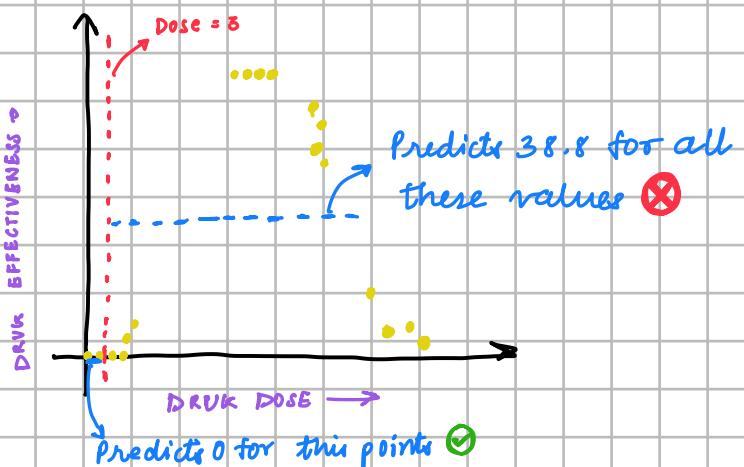
DOSE < 3 ←  
Drug effectiveness = 0  
(avg is zero this side)



Tree will be formed like :-

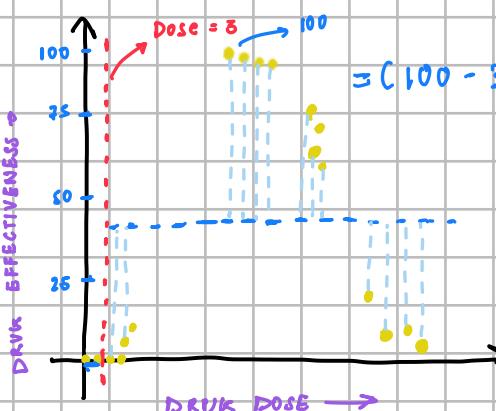


Let's see visual representation of this problem.



From visualization, we can say that this decision tree is not doing good job.

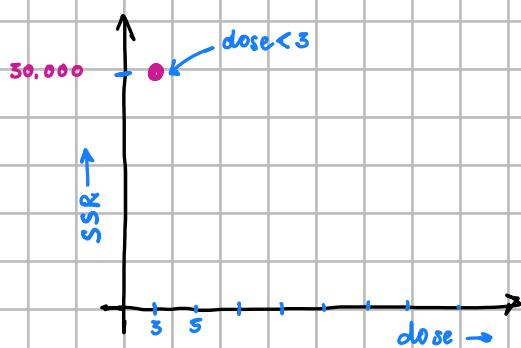
**Step ③ : Sum of squared residuals (SSR)**  
Quantify how good or bad decision tree is!



Residual is a diff. b/w true value & predicted value.

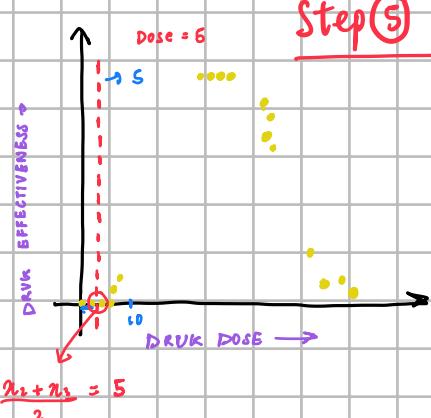
$$\begin{aligned} \text{SSR} = & (0 - 0)^2 + (0 - 38.8)^2 + (0 - 38.8)^2 \\ & + (5 - 38.8)^2 + (20 - 38.8)^2 + \\ & (100 - 38.8)^2 + \dots + (0 - 38.8)^2 \\ & = 27468.5 \end{aligned}$$

**Step ④ : Plot SSR graph with dose on x axis & SSR on y axis**



But, we gotta know that dose < 3 & this is not our root node.

We will repeat this step of SSR until we get accurate value.



**Step ⑤ : Repeat the above for different threshold.**

for dose < 5  
drug effectiveness = 0  
(again both points is at zero value)

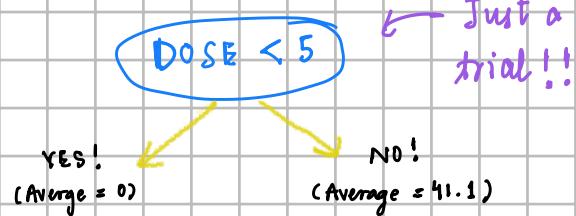
for dose > 5  
drug effectiveness = Avg for all points on right side.

$$= 41.4$$

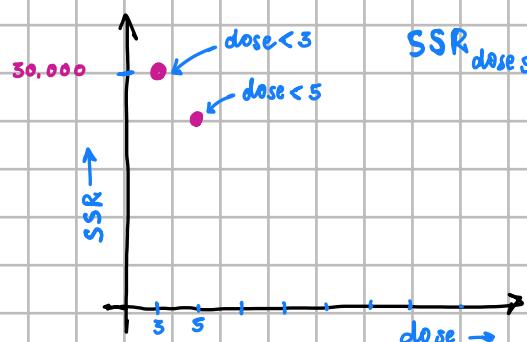
$$\frac{n_1 + n_2}{2} = 5$$

# Visualization !

Tree will be formed like :-



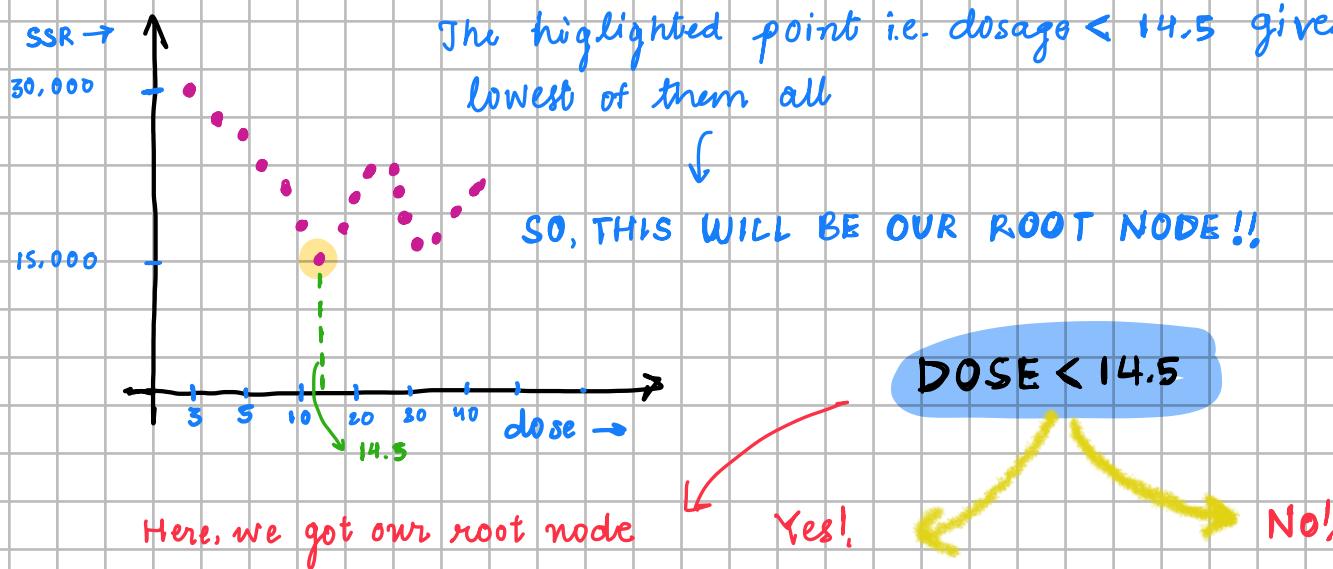
Let's plot SSR graph for s as well.



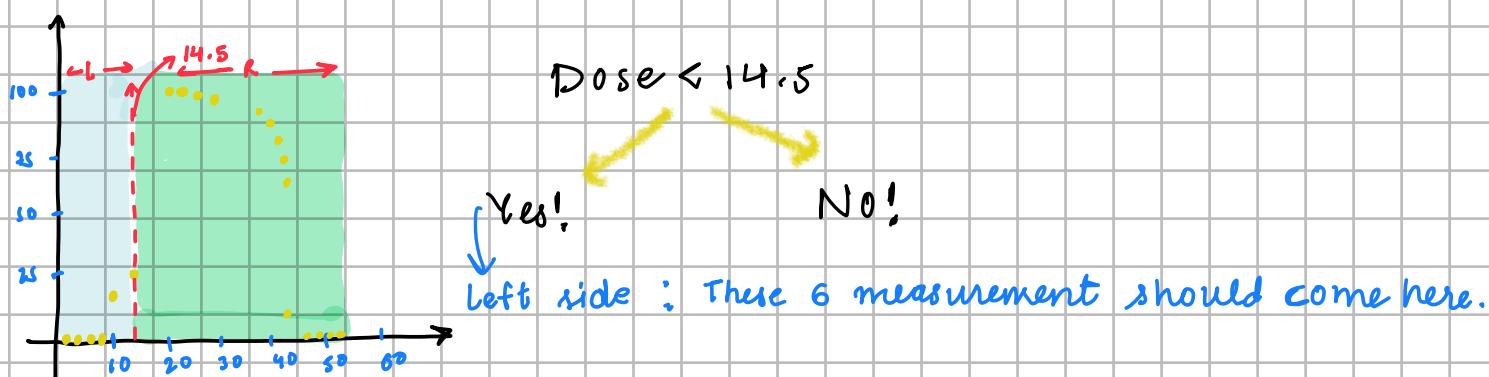
Calculate for different doses & plot in this same graph to get value of root node.

⑥ similarly, calculate for dose < 7, dose < 9, dose < 14.5, dose < 26 and dose < 36

⑦ Plot them altogether in SSR graph :-

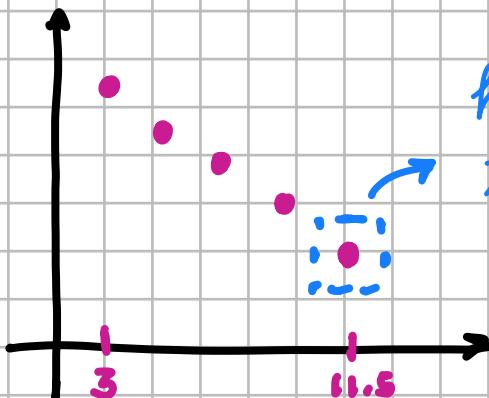


Step ⑧ : Our next goal is to find internal nodes

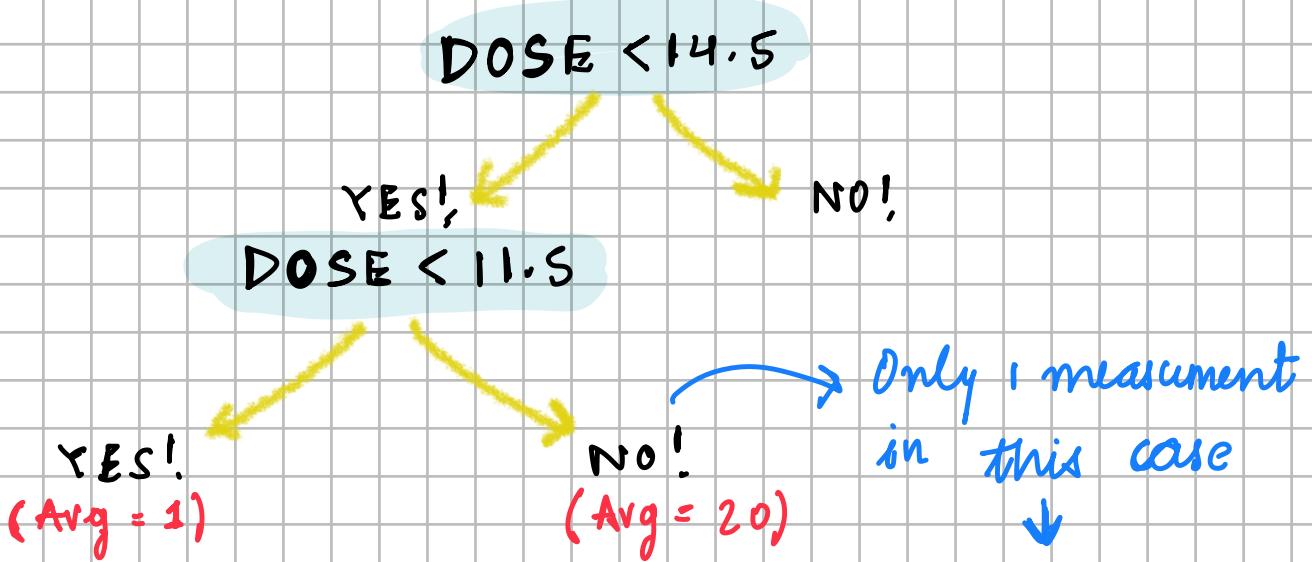


Step 9: Repeat SSR steps we did earlier for these 6 measurements.

Step 10: Plot the graph



here, dose  $< 11.5$  gives lowest threshold  
(previously done in 14.5)



Overfitting!

Step 11: Preventing overfitting

To prevent overfitting, we have to set the minimum no. of measurements for splitting tree.

Let's choose no. as 7

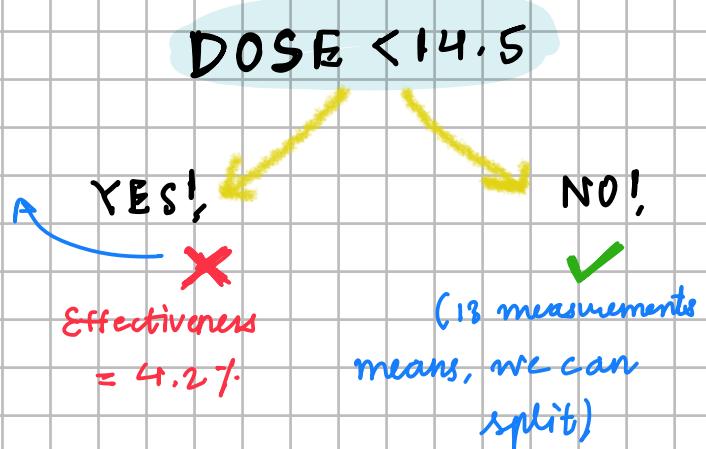
As, we know there are only 6 measurements

(we cannot split it further)

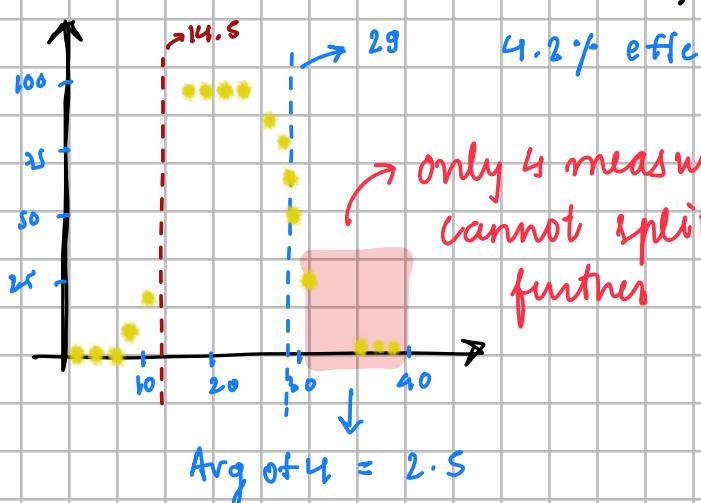
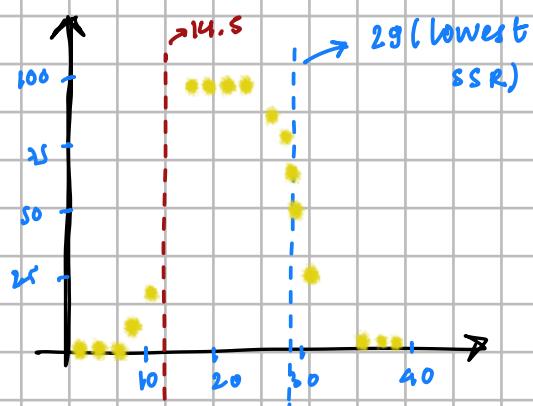
Effectiveness = avg of 6 measures.

We will split the measurements from right side.

(into two groups & find lowest SSR)



Again finding SSR values :



By this, we gotta know that our next internal node (right) would be dose  $> 2.9$

$\text{DOSE} < 14.5$

YES!

NO!

$\text{DOSE} > 2.9$

YES!

$2.5\%$  effective

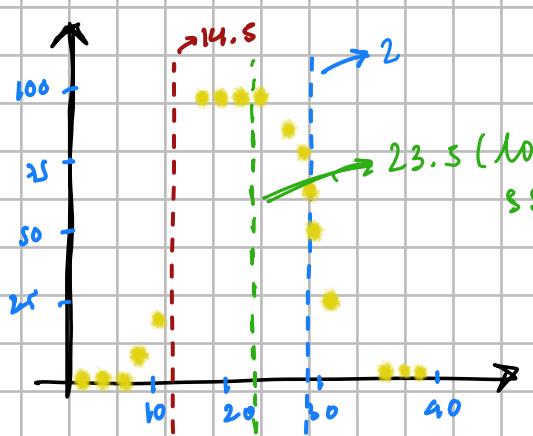
NO!

hot's chuk!

(7 measurements)

\* Now, again we'll have to decide measurement to split b/w 14.5 and 2.9.

\* Following same procedure, we will find lowest SSR.



$\text{DOSE} < 14.5$

YES!  
 $4.2\%$  effective

NO!  
 $\text{DOSE} > 2.9$

YES!  
 $2.5\%$  effective

NO!  
 $\text{DOSE} > 23.5$

\* Here, both the values have lower value than 7. So we cannot split this further

Our final tree structure is ready !!

YES!  
 $52.5\%$  effective  
NO!  
 $100\%$  effective

# Multiple features in this tree

Using dose, age and sex-

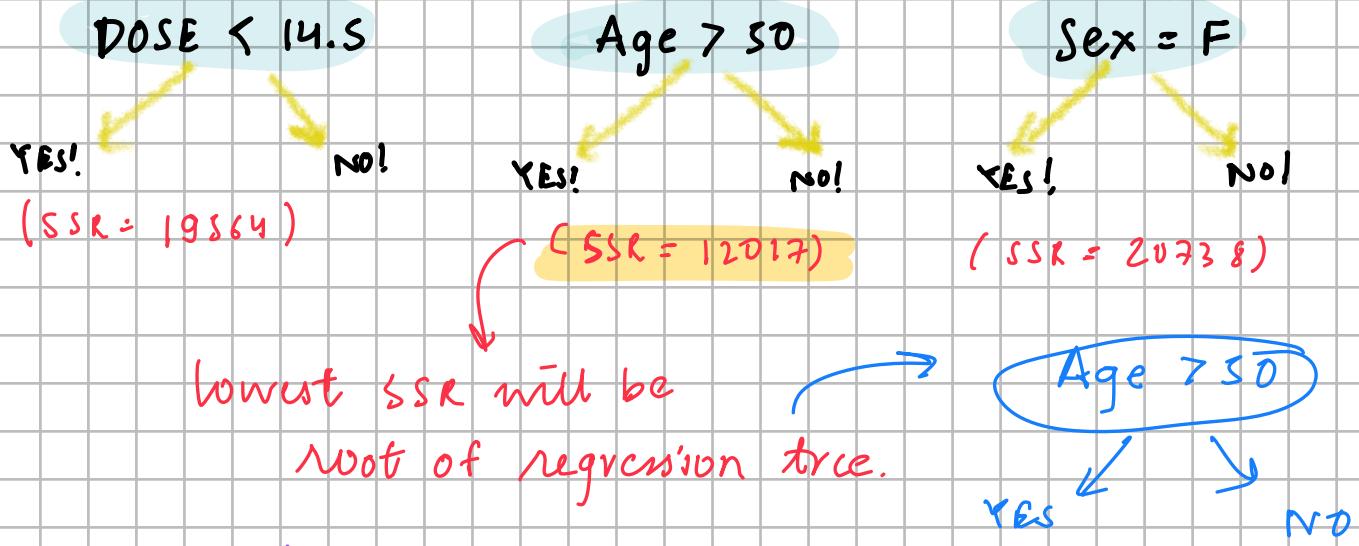
- \* First, we'll completely ignore age and sex and use dose :-

Dose	Age	Sex	Drug Effectiveness
10	25	F	98
20	73	M	0
35	54	F	6
5	12	M	44

- \* We select threshold with smallest sum of squared residuals (SSR) →

- \* Then, we ignore dose and sex, and only age to predict drug effectiveness.

- \* By this, we individually find SSR and create preferences according to it.



- \* Then we grow the tree just like before the (only diff is that we'll have 3 candidates)

Do this until we can no longer split further.

