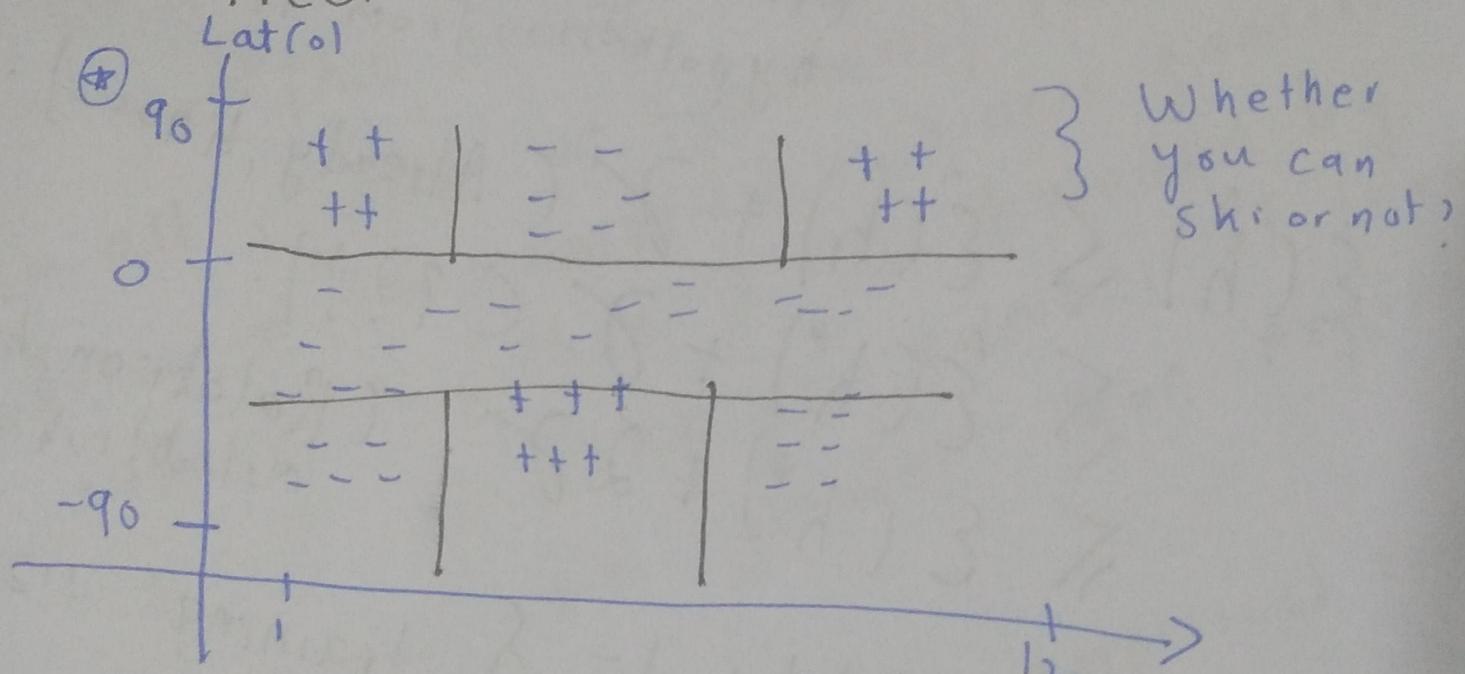


## Lecture-10

### ④ Decisions Tree:-

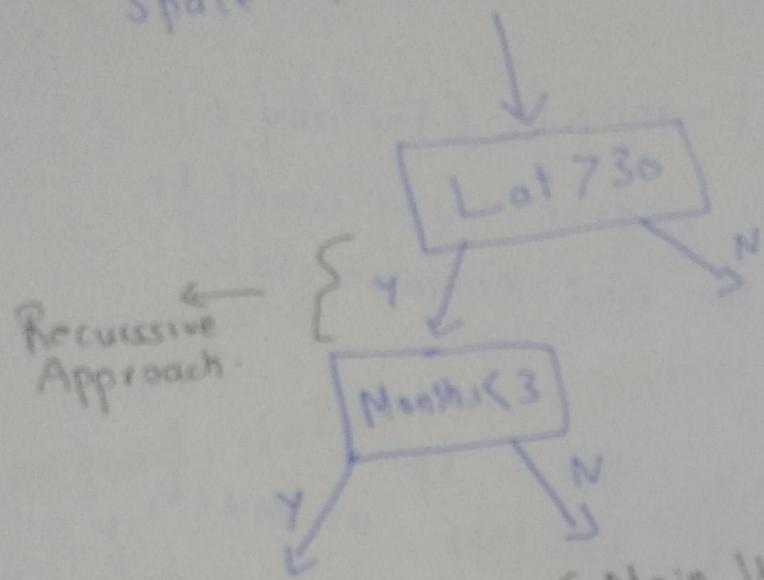


- ④ Difficult to come up with linear decision boundary that is optimal
- ④ Decision tree is optimal in this case.

(112)

→ Decision Tree works by greedy,  
Top-Down - Recursive Partitioning

→ Decision is gonna ask question in the  
 Space to create isolated regions (Small  
 spatial regions).



→ More formally: (Main Working)  
 Region  $R_p$

Looking for a split  $S_p$

$$S_p(j, t) = (\{x_i | x_{ij} < t, x \in R_p\}) \quad \text{①}$$

$$= (\{x_i | x_{ij} \geq t, x \in R_p\}) \quad \text{②}$$

Feature  
to be  
used  
for split

Threshold

① How to choose splits:-

Define  $L(R)$ : Loss on R

Given  $C$  classes, define  $\hat{P}_c$  to be the proportion of examples in R that are of class c.

②  $L_{\text{misclass}} = 1 - \max_c \hat{P}_c$  (For any sub-divided region you want to pick up the most common class there).

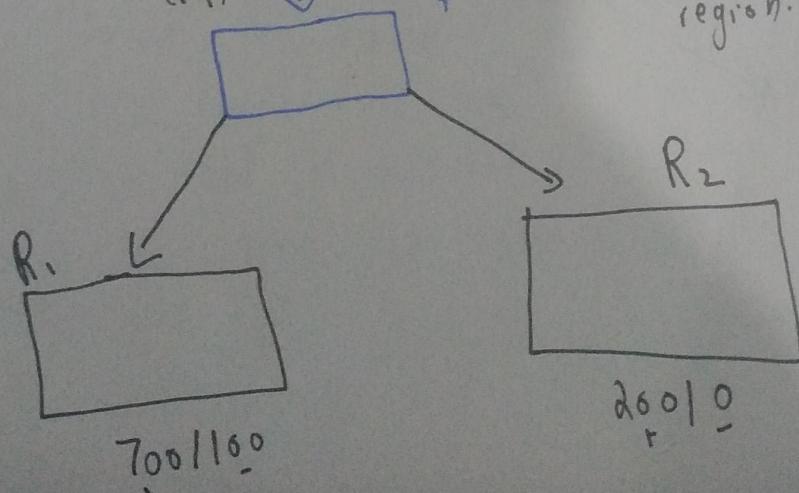
③ Want to pick the split which decreases the loss as much as possible:

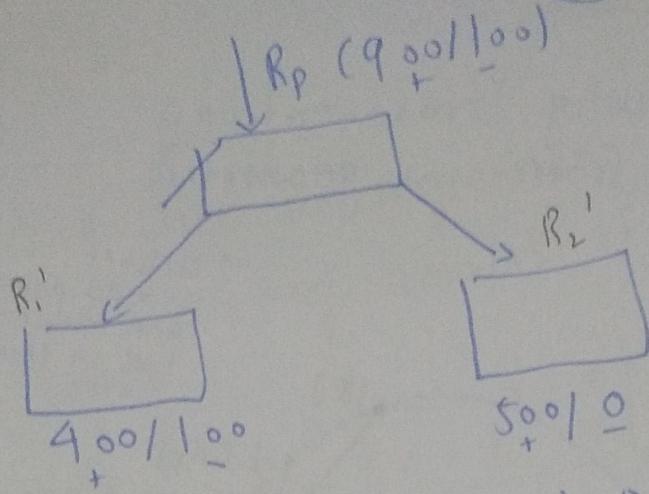
$$\max_{R_1, R_2} L(R_p) - (L(R_1) + L(R_2))$$

↓ parent loss      ↓ children loss

Loss has issues:-

(R<sub>p</sub>)  $\downarrow$   $900/100 \rightarrow$  Number of examples in this region.





$$\begin{aligned} L(R_1) + L(R_2) &= 100 + 0 = 100 \\ L(R_1') + L(R_2') &= 100 + 0 = 100 \end{aligned} \quad \rightarrow \begin{array}{l} \text{Loss} \\ \text{Remains} \\ \text{the same} \\ \text{in above} \\ \text{two splits.} \end{array}$$

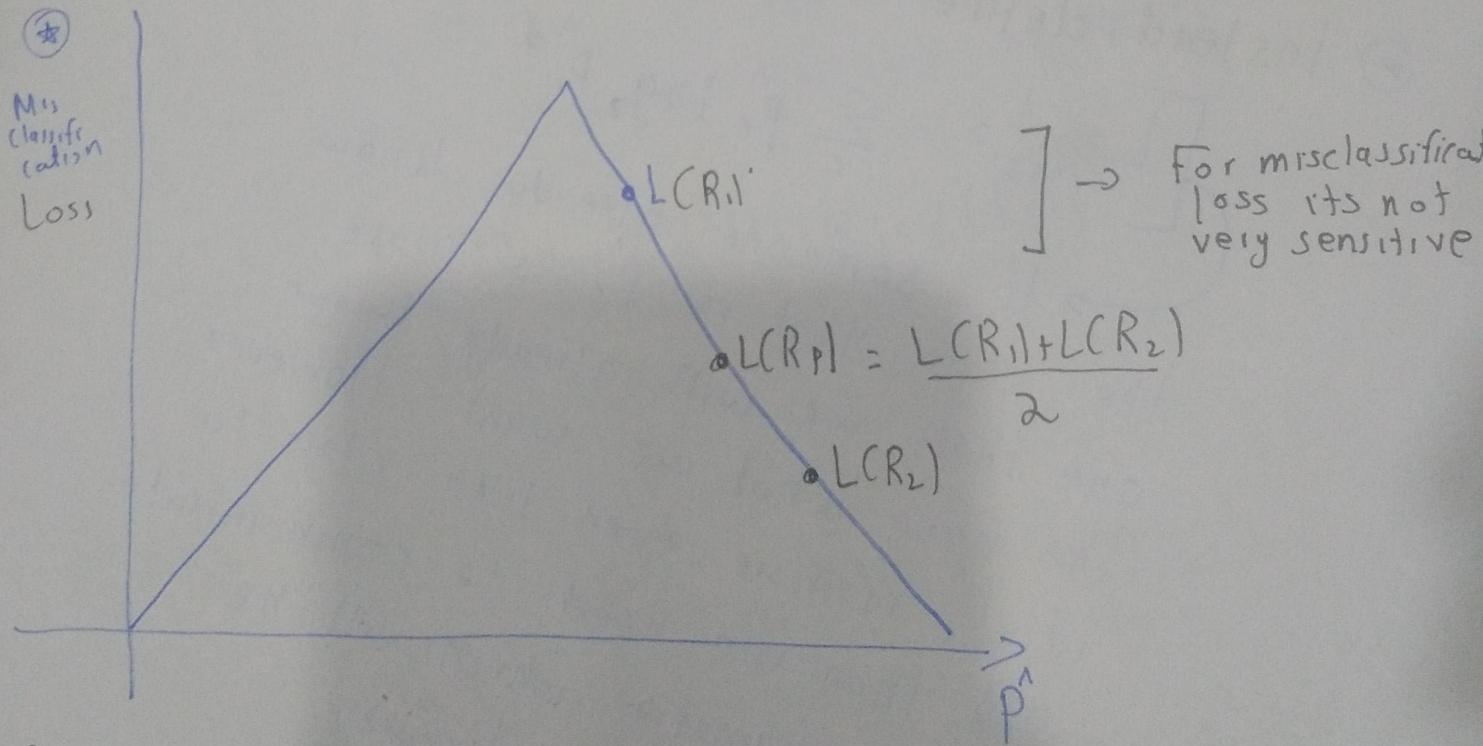
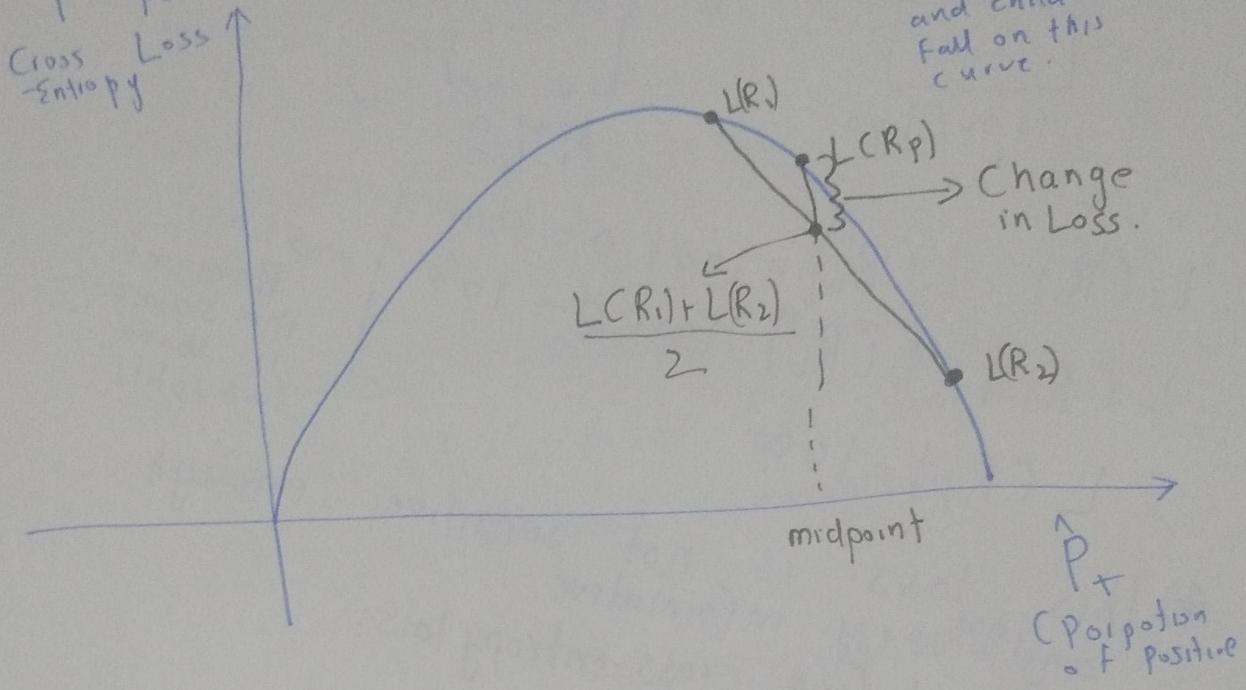
$$L(R_p) = 100.$$

- ⑤ This loss is not sensitive enough  
and very informative
- ⑥ Instead, define cross-entropy loss:

$$L_{\text{cross}} = - \sum_c \hat{P}_c \log_2 P_c$$

→ If someone already know  
the class that is  $\hat{P}_c$  is close  
one that we will already know  
if  $\hat{P}_c$  is more even like close to  
0.5 than we need to communicate  
more information which will be  
visible by larger  $L_{\text{cross}}$ .

④ We can see why cross-entropy loss is better from a geometrical perspective:-

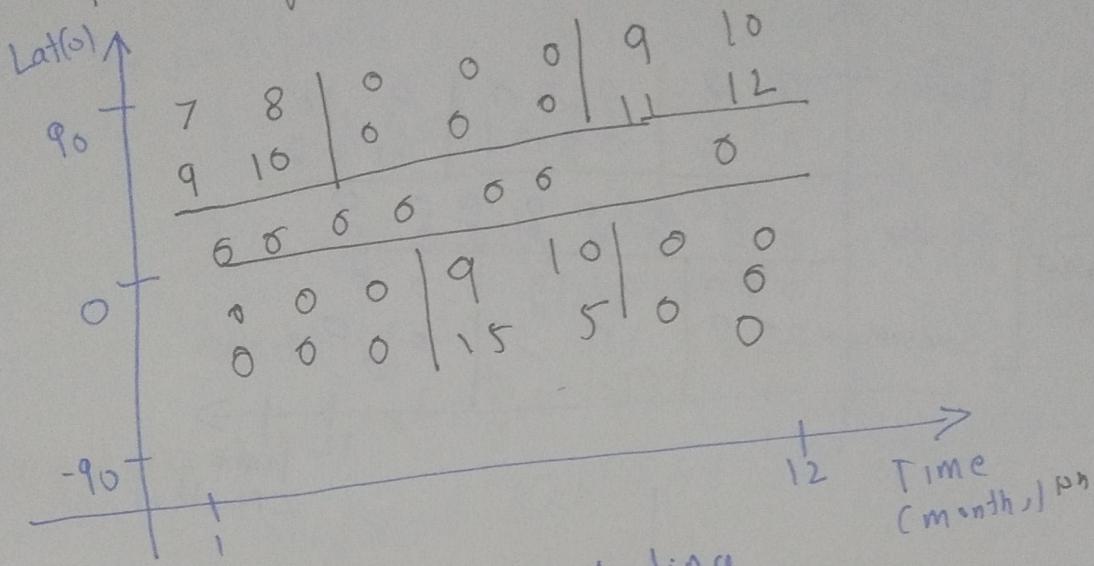


④ Gini Loss:-

$$\sum_c \hat{P}_c (1 - \hat{P}_c) \rightarrow \text{Curve is similar to cross-entropy loss}$$

(116)

④ Decision Tree for regression:  
 → Predicting amount of snowfall:-  
 (in inches).



↳ Instead of predicting the majority of values → You can predict the mean of values present in the region.

$\hat{y}_m$  (Mean is the prediction)

$$\hat{y}_m = \frac{\sum_{i \in R_m} (y_i - \hat{y}_m)^2}{|R_m|}$$

$L_{\text{square}} = \frac{\sum_{i \in R_m} (y_i - \hat{y}_m)^2}{|R_m|}$  [C Loss]  
 ↓  
 of specific region.

(115)

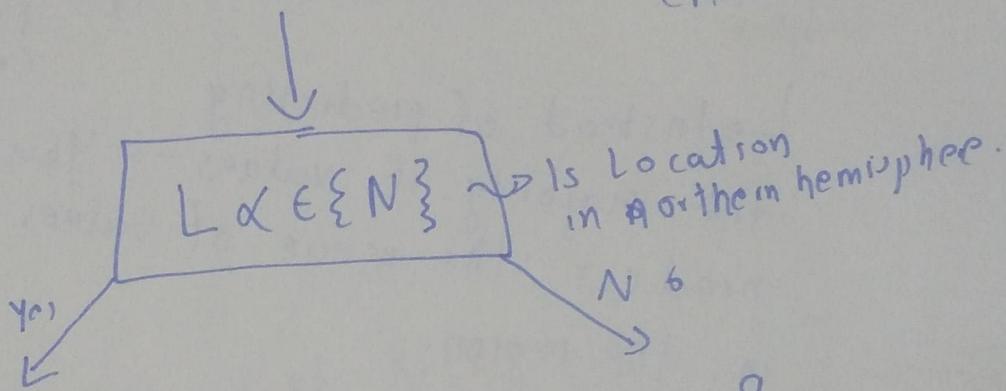
(117)

\*)

Categorical Variables can be dealt  
efficiently:-

	Northern			Equator			Southern			
	7	8	9	0	0	0	9	10	11	12
	9	10		0	0	0				
	0	6	0	6	6	6	0	0	0	0
	0	6	0	9	10		0	0	0	0
	0	6	0	15	5		0	0	0	0
			!							12
										Time (Months)

Categorical  
Variable.



→ However, remember if you have  $q$  categories than you need to consider  $2^q$  splits which is computationally efficient.

⊕ If you let your tree grow indefinitely then you would have separate region for each point

↳ Probably you are overfitting

↳ Generally decision trees have tendency for high variance.

tendency for high variance.  
of Decision Trees:-

• Regularization of Decision Trees → You stop

1) Minimum Leaf size → You stop splitting

2) Max depth

3) Max number of nodes

4) Minimum decrease in Loss  
↳ You might stop

too early.

5) Grow the full Decision

tree and then prune it

backwards. (By using its misclassification with validation set).

## ④ Runtime:

- n examples
- f features
- d depth

Test time: → pretty quick  
 $O(d)$  → Runtime only  
your depth; typically  
 $d < \log n$

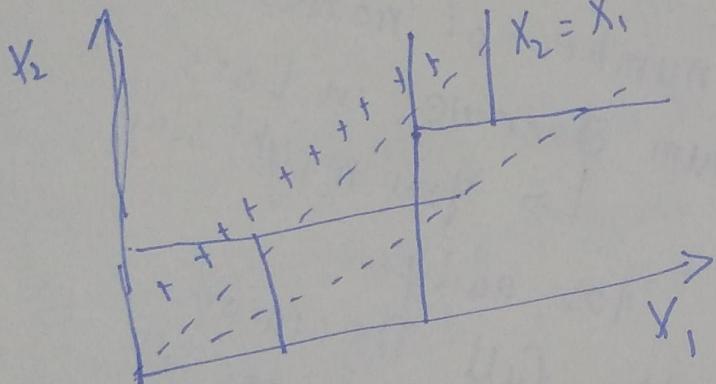
## Train Time:

→ Each point is part of  $O(d)$  nodes  
(Cost of point at each node is  $O(f)$ .)

Total Cost:  $O(nfd)$

↳ Not that slow

• Decision trees don't have additive structure



• Decision tree  
have problem  
of over  
complicating  
when feature(s)  
follow an  
additive  
structure.

## ⑤ Recap:

### ⊕ Benefits:

⊕ Easy to explain

⊕ Interpretable

⊕ Can deal with categorical variable

⊕ Fast

### ⊖ Negative:

⊖ High Variance    ⊖ Bad at additive

⊖ Low predictive accuracy

- ④ You can make decision tree stronger through **ensembling**

### → Ensembling:

→ Take  $X_i$ 's which are random variable (RV) that are independent identically distributed (IID)

$$\rightarrow \text{Var}(X_i) = \sigma^2, \text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum X_i\right)$$

$\downarrow$

$= \frac{\sigma^2}{n}$

Mean  
of many  
of these  
variables

$\downarrow$   
Increase  
in  $\bar{Y}_i$ 's  
variance

→ What if you drop the independent assumption?

So now  $X_i$ 's are iid (identically distributed)

$X_i$ 's correlated by  $P$

$$\text{Var}(\bar{X}) = P \sigma^2 + \frac{1-P}{n} \sigma^2$$

→ This shows we want more examples but make sure that as many examples are decorrelated.

## ① Different Ways to ensemble:-

- 1) Different Algorithm
- 2) Different Training Sets
- 3) Bagging  $\xrightarrow{\text{Random Forest}}$  Our Focus
- 4) Boosting  $\xrightarrow{\text{Ad Boost, Xg boost}}$

## ② Bagging:-

### ③ Bootstrap aggregation

↳ Typical used in statistics  
to measure the uncertainty of  
your estimate.

→ What is Bootstrap aggregation?

→ Have a true population  $P$

→ Training set  $S \sim P$  ( $S$  sampled from  $P$ )

→ Assume your population is  
training sample ( $S = P$ ).

→  $\hookrightarrow$  Bootstrap Samples  $Z \sim S$  :-

↳ You sample  $N$  example  
from your training set  $S$  with  
replacement.

→ Train your models on  $Z$  and  
see the variability that model is making  
based on different bootstrap examples.

→ We are going to bootstrap sample, train different models and average their output (Bootstrap aggregation).

→ Let's make it more formal:

① Bootstrap Samples  $Z_1, \dots, Z_M$

② Train model  $G_m$  on  $Z_m$ .

$$G(m) = \underbrace{\sum_m G_m(x)}_{M} \quad \begin{array}{l} \rightarrow \text{Aggregating} \\ \text{output of} \\ \text{all models.} \end{array}$$

$\rightarrow$  Total Number of Models

→ Bias-Variance Analysis:

Remember

$$\text{Var}(X) = p\sigma^2 + \frac{1-p}{N}\sigma^2$$

- Bootstrapping is driving down p

- Bootstrap decreases M.

→ Less Variance in your model (Less Overfitting).

• But you are increasing the bias of model because you are training on less data (Each model which reduces its complexity).

Decrease in variance is greater than increase in Bias!!!

## ④ DT + Bagging:-

→ Recall:- DT are high variance?  
 low bias → suitable for  
 Bagging (ideal fit).

→ Random Forest  $\rightarrow$  DT + Bagging  
 but it includes more randomization  
 in each tree.

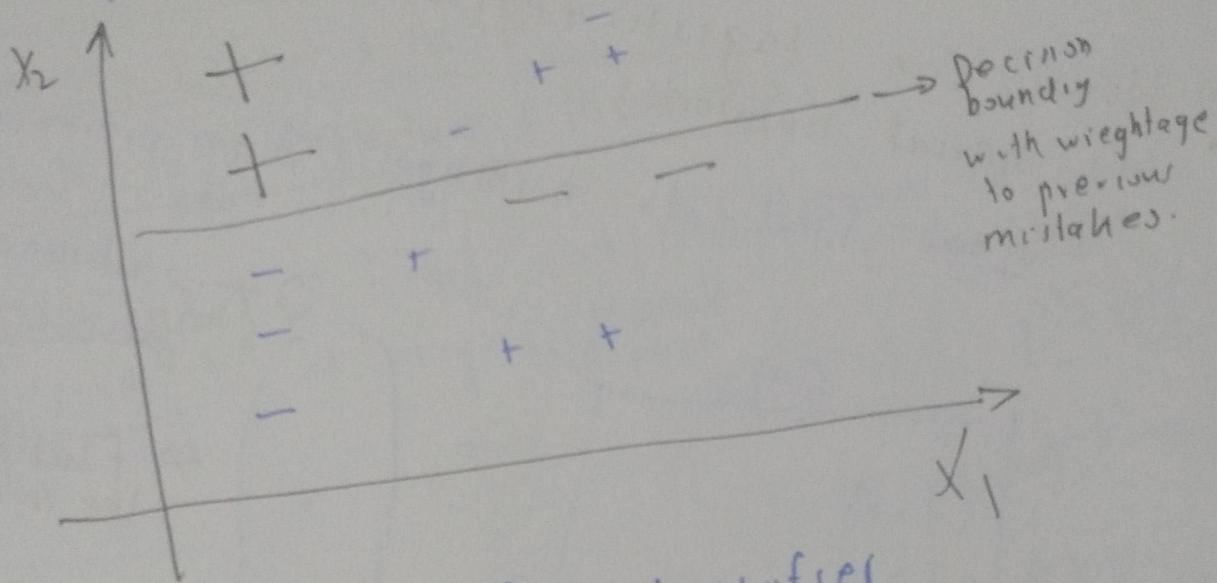
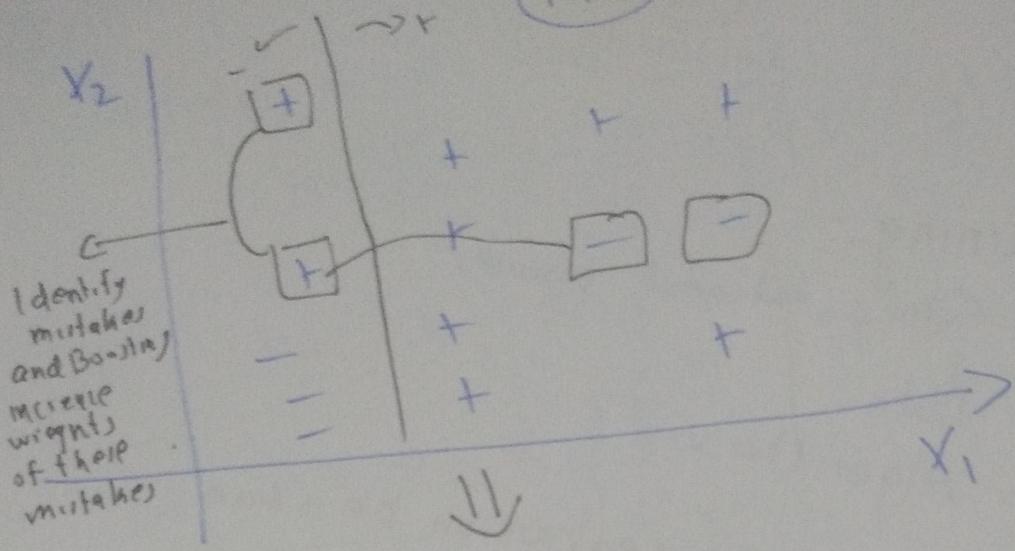
↳ At each split, consider  
 only fraction of your features.  
 $(\text{total features}) \cdot (1 - \text{correlation}$   
 $p)$ . → Make sure different  
 DT in random forest are  
 decorrelated.

## ⑤ Boosting:-

→ Decreasing Bias of your  
 model

→ More additive in nature

→ You train one and you  
 add that prediction into  
 your ensemble.



→ Determine for classifier  
 $\alpha_m$  a weight for each model trained, higher  $\alpha_m$  for good model,  
 lower  $\alpha_m$  for bad model.

↳ In adaboost for example  
 $\alpha_m$  is proportional to:  $\log \frac{1 - \text{err}_m}{\text{err}_m}$

Finally:

$$G_T(x) = \sum_m \alpha_m G_m$$

Each  $G_m$   
 trained on  
 re-weighted  
 training set