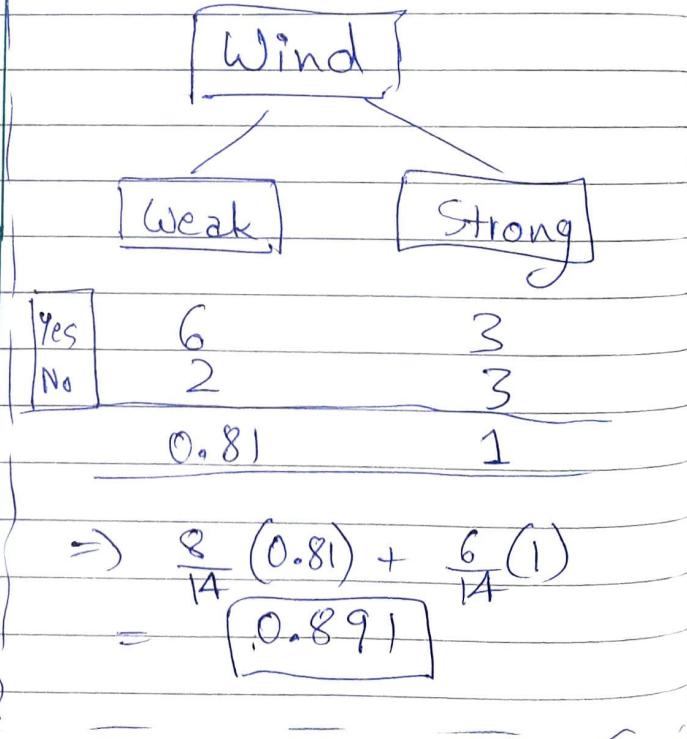
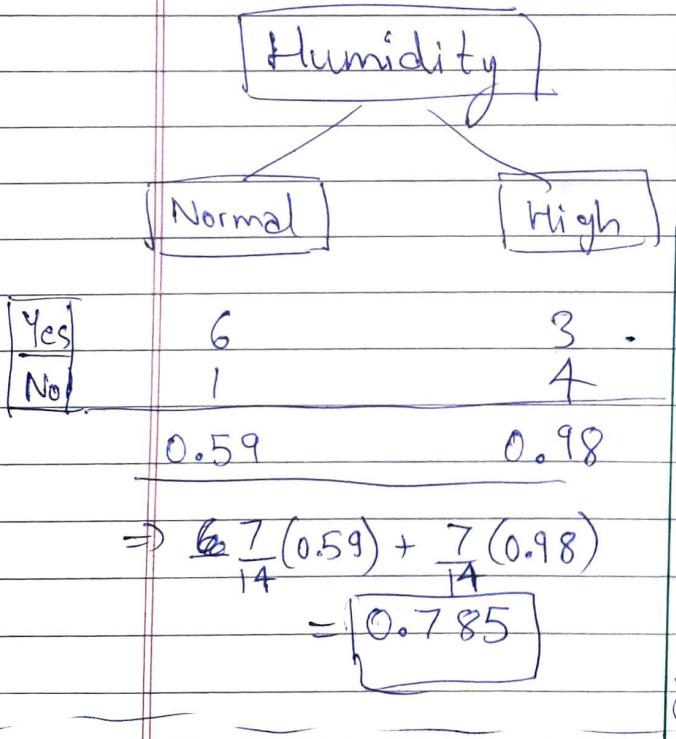
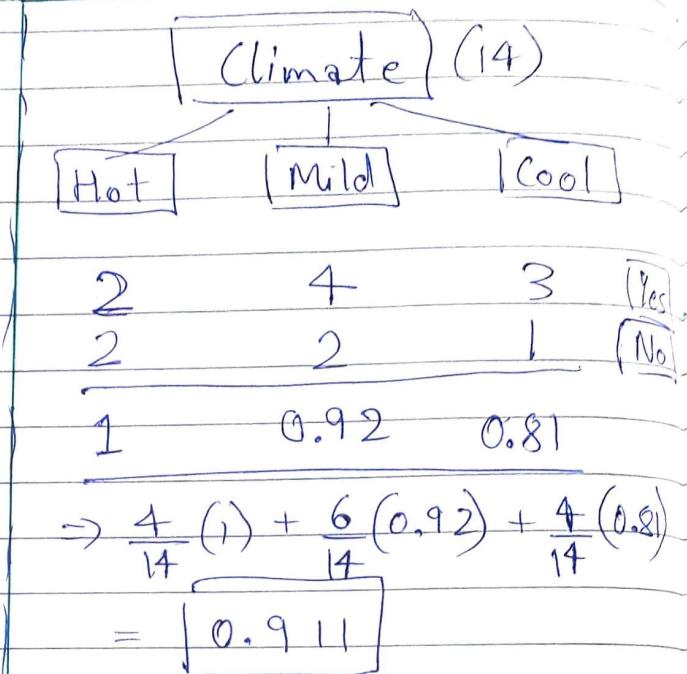
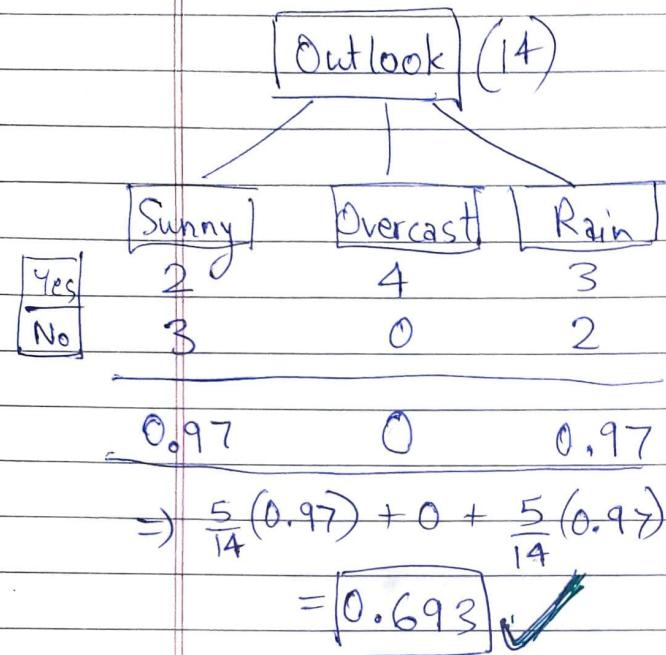
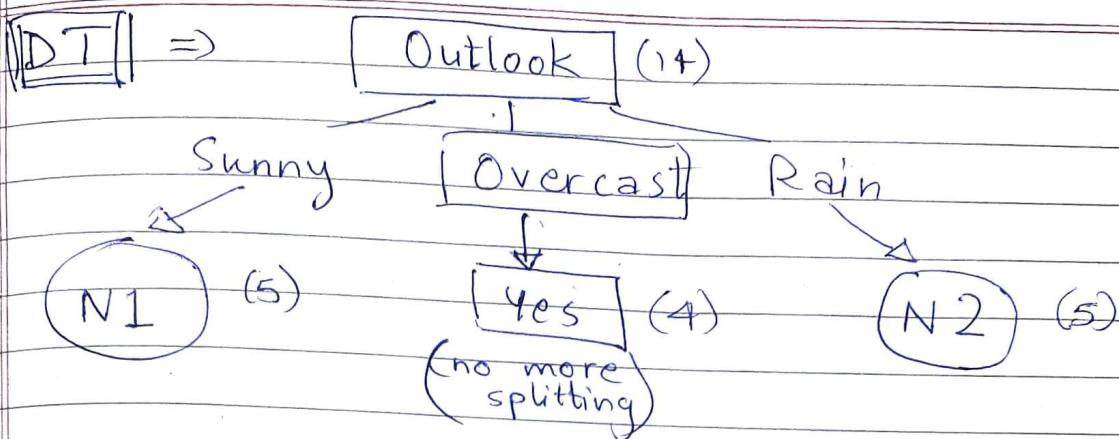


5)(a)

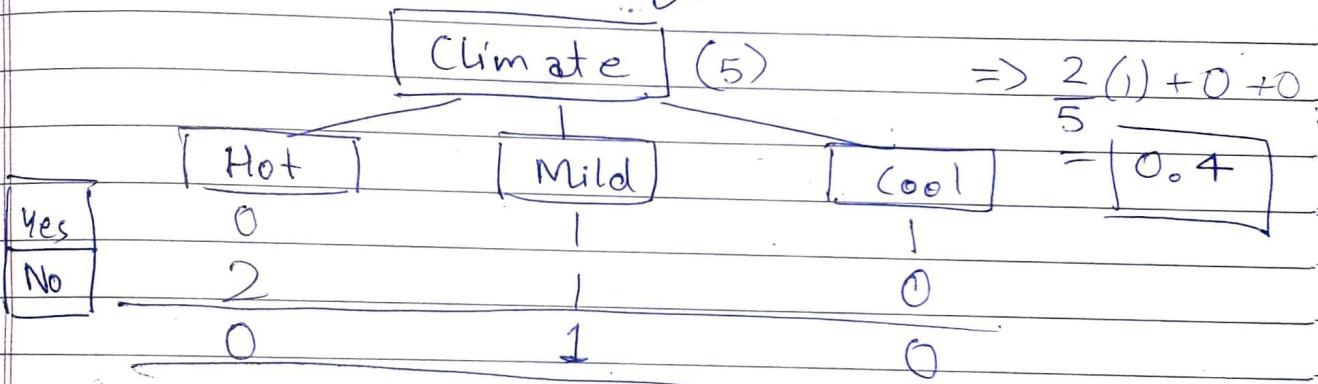
>Selecting the ~~red~~ attribute for split



OUTLOOK GIVES BEST SPLIT



Node N1 (Outlook = sunny)



Humidity (5)

Normal

Yes	2
No	0

0	0
---	---

$\Rightarrow 0$ ✓

Wind (5)

Weak

Yes	1
No	2

0.92

$$\Rightarrow \frac{3}{5}(0.92) + \frac{2}{5}(1) = 0.952$$

DT \Rightarrow

N1 (Sunny)

Humidity = Normal

Yes (2)

Humidity = High

No (3)

▣ N2 (Outlook = Rain)

(Climate) (5)

[Hot] [Mild] [Cool]

$$\Rightarrow 0 + \frac{3}{5}(0.92) + \frac{2}{5}(1)$$

Yes	0
No	1
	1

0	2	1
0	1	1
0	0.92	1

$$= \boxed{0.952}$$

Humidity

[Normal]

[High]

Wind ~~Weak~~

[weak]

[Strong]

Yes

2

1

Yes

0

No

1

1

No

2

$$\Rightarrow \frac{3}{5}(0.92) + \frac{2}{5}(1)$$

$$= \boxed{0.952}$$

Yes

No

0

0

0

0

0

✓



DT

N2 (Outlook = Rain)

Wind=Weak

Yes (3)

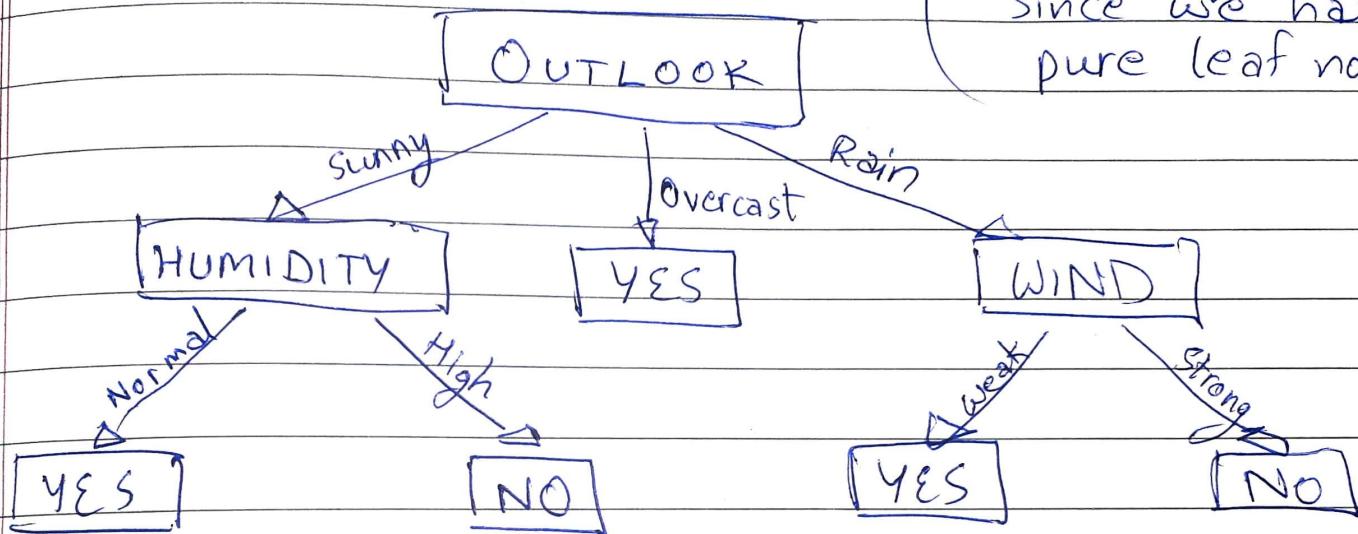
Wind=Strong

No (2)

⇒ FINAL DT will be :-

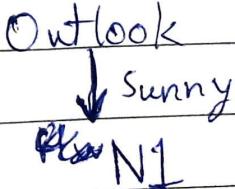
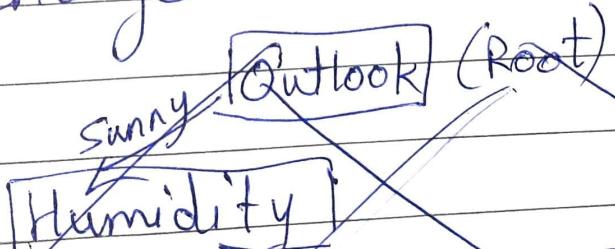
* Accuracy = 100%

Since we have
pure leaf nodes



5b) Outlook = Sunny, Climate = Cool, Humidity = High,
Wind = ~~too~~ Strong, Play = ~~no~~ Yes, Day = D15.

Adding the above row, ~~as~~ the distribution
will change as follows :-

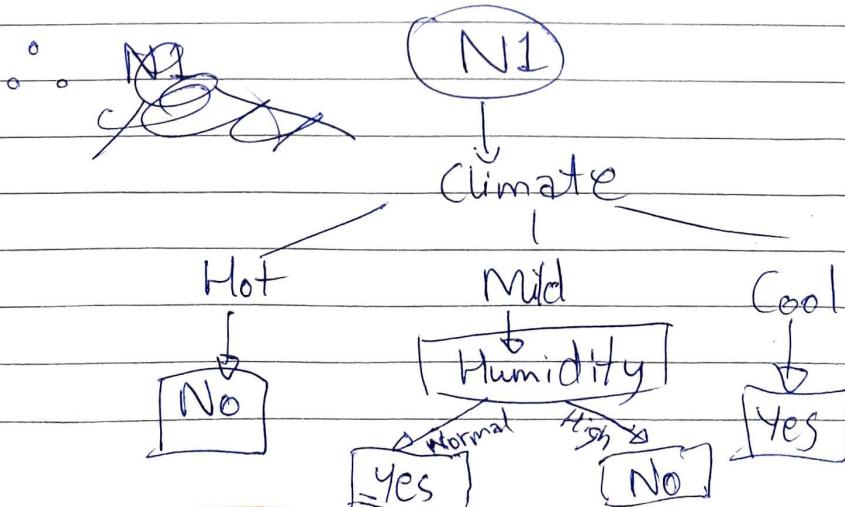


~~Ans~~ Humidity will be selected to split N1. ~~not~~

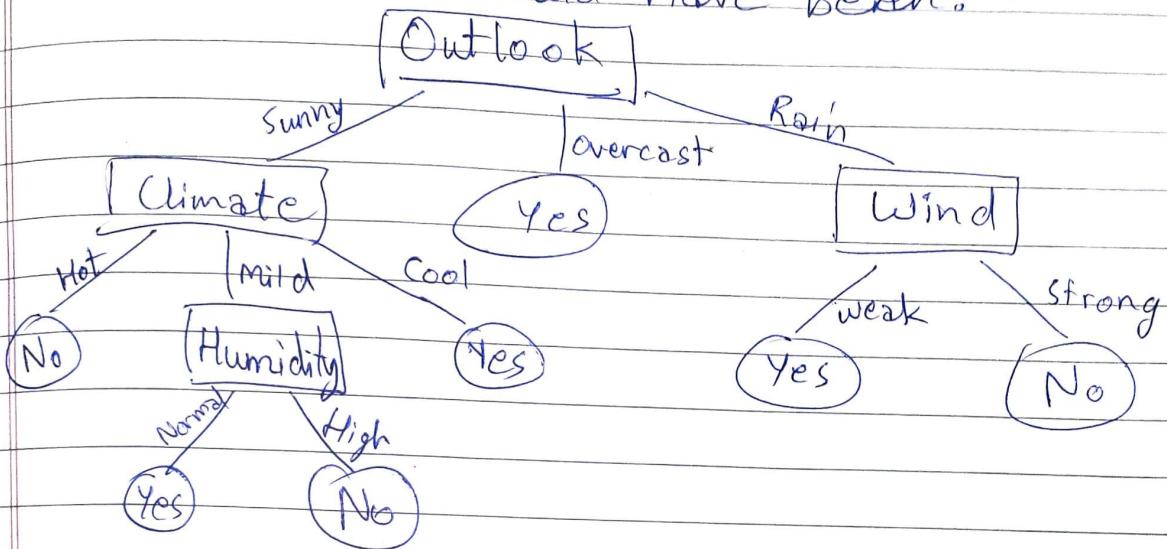
		Humidity		
		Normal	High	$\Rightarrow \frac{4}{6} (.81)$
Yes	2	1		
	No	0	3	$= [0.54]$
		<u>0</u>	0.81	

		Wind		
		Weak	Strong	$\Rightarrow \frac{3}{6} (.92) \times 2$
Yes	1	2		
	No	2	1	$= [0.92]$
		<u>0.92</u>	0.92	

Climate				
			$\Rightarrow \frac{2}{6} (1) = [0.33]$	\checkmark
Hot	Mild	Cool		
Yes	0	1	2	
No	2	1	0	
	<u>0</u>	1	0	



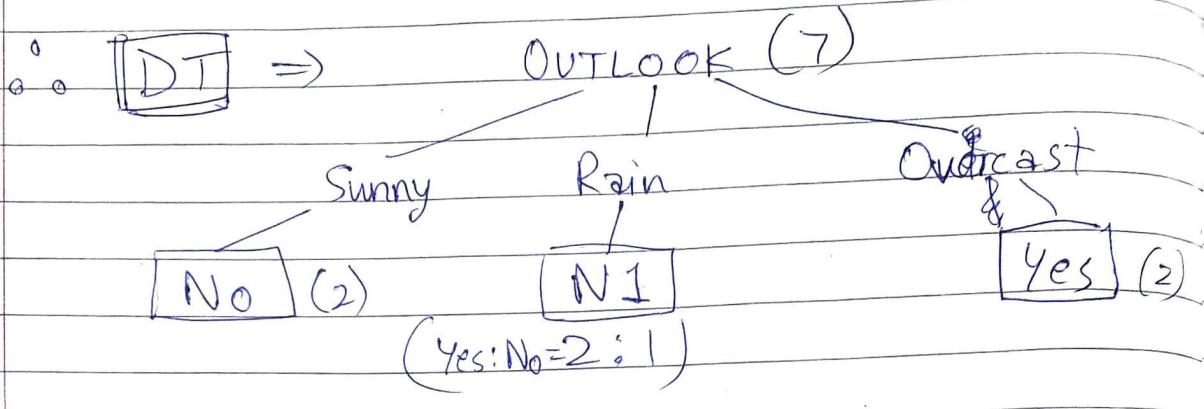
The new DT would have been:-



(5c) D1-D7 (Training)

Outlook (?)			Humidity (?)		
Sunny	Overcast	Rain	Normal	High	
Yes	0	2	2	Yes	2
No	2	0	1	No	1
0	0	0.92		0.92	1
$\Rightarrow (0.92) \frac{3}{7} = [0.394] \checkmark$			$\Rightarrow 0.92 \frac{3}{7} + 1 \frac{4}{7} = [0.966]$		

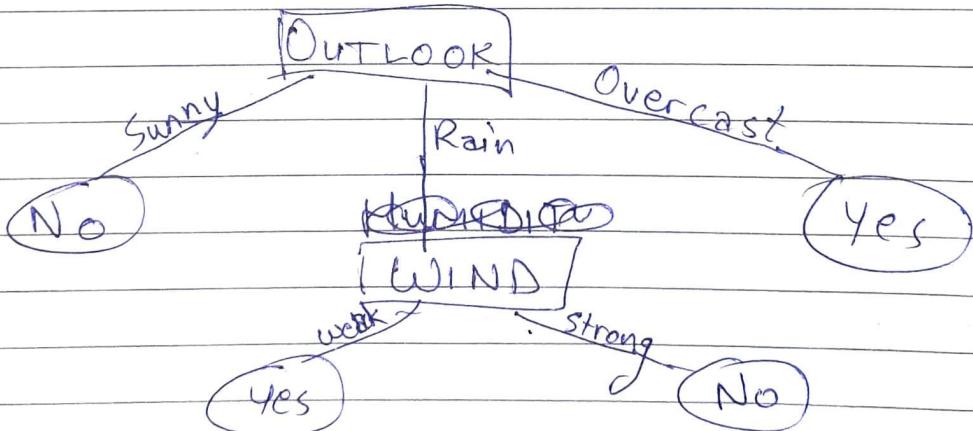
Wind (?)		Climate			
Weak	Strong	Hot	Mild	Cool	
Yes	3	1	1	1	2
No	1	2	2	0	1
0.81	0.92	0.92	0	0.92	
$\Rightarrow \frac{4}{7}(0.81) + \frac{3}{7}(0.92) = [0.857]$		$\Rightarrow \frac{6}{7}(0.92) = [0.788]$			



◻ Outlook = Rain: N1 splitting

By observation, we see that ~~records~~ + records with $\text{outlook} = \text{Rain}$,
 $\text{Wind} = \text{Weak} \Rightarrow \text{Play} = \text{Yes}$
and $\text{Wind} = \text{Strong} \Rightarrow \text{Play} = \text{No}$

∴ DT becomes



* Since leaf nodes are pure; train accuracy = 100%

* Evaluating test accuracy

Node	D8	D9	D10	D11	D12	D13	D14
yPred	No	No	Yes	No	Yes	Yes	No
yTrue	No	Yes	Yes	Yes	Yes	Yes	No

$$\text{Test accuracy} = \frac{5}{7} \times 100 = 71.43\%$$

- * We did not have enough training data, plus we let the tree build fully, thus overfitting on the training data. picking up any noisy sample available also, leading to a bad test accuracy in comparison to the train accuracy.
- * The small amount of data did not allow the model to understand the true relationship.

5(d) There are two possible pruning strategies.

- (i) Pre-pruning :- Stop the algo before fully grown tree, typically if all instances belong to the same class or if all attribute values are the same or by some user specified threshold parameters.
- (ii) Post-pruning:- Generally we do not know when to stop the algorithm. Therefore, first we let the tree grow to its entirety, and then apply pruning.

Since pruning is to avoid overfitting, we should either work with a validation set error or some generalized error measurements.

The following two techniques can be used:-

- Subtree replacement: Trim the nodes in a bottom-up manner, and if the performance improves, replace the subtree with a leaf node, having class label of the majority class
- Subtree raising: Replace the subtree with the branch which is most frequently taken if this improves the ~~generalization~~ performance.

* These strategies would work since the performance measure is taking into account, both seen and unseen data performance.

(6) First ~~the~~ Order Markov model says that Probability of an n^{th} event depends only on the $(n-1)^{th}$ event and the events before that have zero weightage in the probability.

We are supposed to find $P(w_3 | w_1, w_2, w_4)$ which reduces to $P(w_3 | w_2, w_4)$ by the 1st-order Markov assumption.

The order of these events is as follows:-

$$w_2 \rightarrow w_3 \rightarrow w_4$$

\therefore Bayes rule gives

$$P(w_3 | w_2, w_4) = \frac{P(w_4 | w_2, w_3) \cdot P(w_3 | w_2)}{\sum_{\text{possible } w_3} P(w_4 | w_2, w_3) \cdot P(w_3 | w_2)}$$

By Markov assumption,

$$P(w_4 | w_2, w_3) = P(w_4 | w_3)$$

$$\therefore P(w_3 = \text{tough} | w_2, w_4) \quad (\begin{matrix} w_2 = \text{course} \\ w_4 = \text{course} \end{matrix})$$

$$= P(\cancel{w_4} | w_3 = \text{tough}) \cdot P(w_3 = \text{tough} | w_2)$$

$$P(w_4 | w_3 = \text{tough}) \cdot P(w_3 = \text{tough} | w_2) + P(w_4 | w_3 = \text{course}) \cdot P(w_3 = \text{course} | w_2)$$

$$= \frac{(0.3)(0.5)}{(0.3)(0.5) + (0.5)(0.5)} = \frac{3}{8} = 0.375$$

$$\boxed{P(w_3 = \text{tough} | w_2, w_4) = 0.375}$$

$$\boxed{P(w_3 = \text{course} | w_2, w_4) = 1 - 0.375 = 0.625}$$

(7) (a) DT : Decision Tree ; LR : Logistic Regression

(a) \rightarrow Unlike LR, DT does not assume all the features to be independent. ~~This is an~~ ~~drawback~~

\rightarrow But, the biggest advantage is that once built, the DT is easy to understand and gives direct pattern for data analysis. It can be easily converted to an SQL syntax to apply on big datasets.

\rightarrow DT works better ~~when there are more than a single decision boundary, but parallel to the axes.~~ ^{feature}

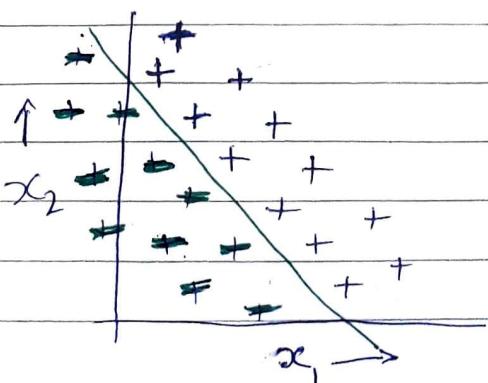
7(b) \rightarrow LR, unlike DT, can draw any decision boundary, which need not be parallel to the ^{feature} axes.

\rightarrow ~~DT~~ DT has relatively high cost in terms of memory and time. Splitting a node involves lots of comparisons, that is kind of a brute force.

\rightarrow DT when scaled up to very high depths are very much prone to overfitting.

$$7(c) \quad y = \begin{cases} +1 & w_1x_1 + w_2x_2 + b \geq 0 \\ -1 & w_1x_1 + w_2x_2 + b \leq 0 \end{cases}$$

7(c) ~~(d)~~ n :- no. of data samples



In the worst case, each data sample gets classified with the help of ~~2~~ unique decision boundaries.

In the worst case, thus, we will have $n+1$ boundaries on each of the axes.

~~This, there will be $(n+1)^2$ leaf nodes.~~

- * Imagine first of all x_1 only is used for creating boundaries leading to $(n+1)$ leaves.
- * Each of these $(n+1)$ leaf nodes would further split to create $(n+1) \times (n+1)$ leaf nodes in total.

Thus, a bound on depth can be found.

Let ~~depth~~ level 0 be the root.

Let level i have 2^i nodes considering only binary splits.

$$\begin{aligned} \therefore 2^i &= (n+1)^2 \\ \Rightarrow i &= \lceil \log_2(n+1) \rceil \leq \log_2(2n) \\ (\text{Depth}) &= \lceil O(\log_2(n)) \rceil \end{aligned}$$

7(d) ~~Answer~~ Continuing from (e) part, if ~~this~~ x_2 were linearly dependent on x_1 , then the boundary for x_2 could be derived from x_1 .

* After $(n+1)$ boundaries by x_1 , we would only need 1 additional node ~~per~~ per leaf to create the x_2 boundary.

* Thus $2(n+1)$ leaf nodes.

$$\begin{aligned} \text{Level } i : \quad 2^i &= 2(n+1) \\ \Rightarrow i-1 &= \lceil \log_2(n+1) + 1 \rceil \leq \log_2(2n) \\ \therefore \text{Depth} &= O(\log(n)) \end{aligned}$$

(8) $X = \langle X_1, X_2, X_3, \dots, X_n \rangle$ $X_i \in \text{Boolean Var.}$
 $P(Y=1) = \prod_{i=1}^n$

To.P.: $P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$

PF: Let $\theta_{ij} = P(X_i=1|Y=1)$ ($\forall 1 \leq i \leq n$)
 $\Rightarrow P(X_i=0|Y=1) = 1 - \theta_{ij}$ ($\forall 1 \leq i \leq n$)

Similarly, $P(X_i=1|Y=0) = \theta_{io}$ ("")
and, $P(X_i=0|Y=0) = 1 - \theta_{io}$ ("")

From the naive bayesian model, we have,

$$\begin{aligned} P(Y=1|X) &= \frac{P(X|Y=1) \cdot P(Y=1)}{\sum_c P(X|Y=c) \cdot P(Y=c)} \\ &= \frac{1}{\frac{P(X|Y=0) \cdot P(Y=0)}{P(X|Y=1) \cdot P(Y=1)} + \frac{P(X|Y=1) \cdot P(Y=1)}{P(X|Y=1) \cdot P(Y=1)}} \\ &= \frac{1}{1 + \frac{1-\pi}{\pi} \cdot \frac{P(X|Y=0)}{P(X|Y=1)}} \end{aligned}$$

* $P(X_i|Y=c)$

$\cancel{X_i=0}$	$\cancel{X_i=1}$
$(1-\theta_{ic})$	(θ_{ic})
$\boxed{(1-\theta_{ic})^{1-X_i} (\theta_{ic})^{X_i}}$	

$$\begin{aligned} &= \frac{1}{1 + \exp(\ln(\frac{1-\pi}{\pi} \cdot \frac{P(X|Y=0)}{P(X|Y=1)}))} \\ &= \frac{1}{1 + \exp(z)} \quad (\text{say}) \end{aligned}$$

$$\begin{aligned}
 z &= \ln\left(\frac{1-\pi}{\pi}\right) + \ln\left(\frac{P(X|Y=0)}{P(X|Y=1)}\right) \\
 &= \ln\left(\frac{1-\pi}{\pi}\right) + \ln\left(\frac{\prod_{i=1}^n P(X_i^0|Y=0)}{\prod_{i=1}^n P(X_i^1|Y=1)}\right) \\
 &= \ln\left(\frac{1-\pi}{\pi}\right) + \sum_{i=1}^n \ln\left[\left(\frac{(1-\theta_{i0})^{1-x_i^0}}{(1-\theta_{i1})^{1-x_i^1}} \frac{\theta_{i0}^{x_i^0}}{\theta_{i1}^{x_i^1}}\right)\right] \\
 &= \ln\left(\frac{1-\pi}{\pi}\right) + \sum_{i=1}^n \ln\left(\frac{1-\theta_{i0}}{1-\theta_{i1}}\right) \\
 &\quad + \sum_{i=1}^n \ln\left[\left(\frac{1-\theta_{i1}}{1-\theta_{i0}}\right) \left(\frac{\theta_{i0}}{\theta_{i1}}\right)\right]^{x_i^0} \\
 &= \underbrace{\ln\left(\frac{1-\pi}{\pi}\right) + \sum_{i=1}^n \ln\left(\frac{1-\theta_{i0}}{1-\theta_{i1}}\right)}_{w_0} + \sum_{i=1}^n x_i^0 \underbrace{\ln\left[\left(\frac{1-\theta_{i1}}{\theta_{i1}}\right) \left(\frac{\theta_{i0}}{1-\theta_{i0}}\right)\right]}_{w_i} \\
 \Rightarrow w_0 &= \ln\left(\frac{1-\pi}{\pi}\right) + \sum_{i=1}^n \ln\left(\frac{1-\theta_{i0}}{1-\theta_{i1}}\right) \\
 \Rightarrow w_i &= \ln\left[\left(\frac{1-\theta_{i1}}{\theta_{i1}}\right) \left(\frac{\theta_{i0}}{1-\theta_{i0}}\right)\right] \\
 &= \ln\left[\frac{P(X_i=0|Y=1)}{P(X_i=1|Y=1)} \cdot \frac{P(X_i=1|Y=0)}{P(X_i=0|Y=0)}\right] \\
 \Rightarrow P(Y=1|X) &= \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i^0)} \quad (+ 1 \leq i \leq n)
 \end{aligned}$$

Hence Proved