

# CS 6375 Machine Learning Midterm Examination

University of Texas at Dallas

10/12/2020

Name: Midterm Solutions

NetID: \_\_\_\_\_

Question	Topic	Points
1	Short Answers	20
2	True/False Questions	8
3	Naive Bayes	10
4	Perceptrons	10
5	SVMs	24
6	Logistic Regression	18
7	Decision Trees	10
<b>Total</b>		100

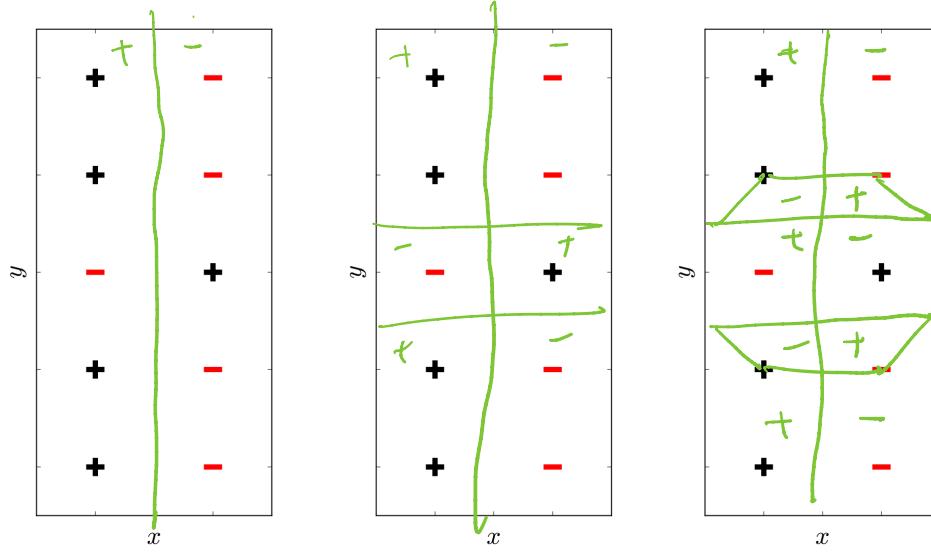
## Instructions:

1. This examination contains 14 pages, including this page.
2. You have **two (2) hours** to complete the examination.
3. Either you can use this paper or separate set of sheets to fill in your answers. Write clearly so we can understand your handwriting. Please scan the answers once you are done (you can also take a photograph once you are done) and upload the scanned copy to eLearning. You can also use a tablet device (like an Ipad) to write if you prefer.
4. Please do not search online for answers to the questions. If the answers are similar to something available online, you will get zero points on this examination.
5. The mid-term examination has to be done individually by everyone. If someone copies, the entire group of students involved will get a zero.
6. Work efficiently. Some questions are easier, some more difficult. Be sure to give yourself time to answer all of the easy ones, and avoid getting bogged down in the more difficult one.
7. All the Best!!

## Question 1: Short Answers

[20 pts] Please provide short and clear answers for the questions below.

- (a) Given a 2 dimensional dataset below, draw the decision boundaries obtained by the following classifiers  
 (4 points: a and b are 1 point each and c is for 2 points)



(a) Logistic regression ( $\lambda = 0$ )

(b) 1-NN

(c) 3-NN

- (b) (6 points: each sub-part is 2 points each) A random variable follows an exponential distribution with parameter  $\lambda : \lambda > 0$ , and has the following density:

$$p(t) = \lambda e^{-\lambda t}, t \in [0, \infty] \quad (1)$$

This distribution models waiting times between events. Given a iid data:  $T = (t_1, \dots, t_n)$ , where each  $t_i$  is modeled as drawn from the exponential distribution with parameter  $\lambda$ . Then:

- Compute the log-likelihood  $p(T|\lambda)$
- Solve for  $\lambda_{MLE}$
- Suppose we have a prior distribution:  $p(\lambda) \propto e^{-\mu\lambda}$ . Obtain  $\lambda_{MAP}$ . Compare  $\lambda_{MLE}$  and  $\lambda_{MAP}$  as  $n \rightarrow \infty$ .

Log-Likelihood

$$p(T|\lambda) = \prod_{i=1}^n p(t_i)$$

$$\begin{aligned}\Rightarrow \log p(T|\lambda) &= \sum_{i=1}^n \log(\lambda e^{-\lambda t_i}) \\ &= \sum_i (\log \lambda - \lambda t_i) \\ &= n \log \lambda - \lambda \sum_{i=1}^n t_i\end{aligned}$$

$$\frac{d \log p(T|\lambda)}{d \lambda} = \frac{n}{\lambda} - \sum_{i=1}^n t_i \geq 0$$

$$\Rightarrow \lambda_{MAP} = \overline{\sum_{i=1}^n t_i}$$

$$\overbrace{p(\lambda) \propto e^{-\lambda}}^{\text{MAP Estm}}$$

$$\begin{aligned}\Rightarrow \log p(\lambda|T) &= -n\lambda \\ &\quad + \sum_{i=1}^n \log \lambda e^{-\lambda t_i} \\ &= n \log \lambda - \lambda \sum_{i=1}^n t_i - n\lambda \\ &= 0\end{aligned}$$

$$\Rightarrow \lambda_{MAP} = \overline{\sum_{i=1}^n t_i + n\lambda}$$

as  $n \rightarrow \infty$

$$\lambda_{MAP} \rightarrow \overline{T} M V \bar{\lambda}$$

(c) (6 points) Decision Trees: Using the dataset below, we want to build a decision tree which classifies Y as T or F given the binary variables A, B, C.

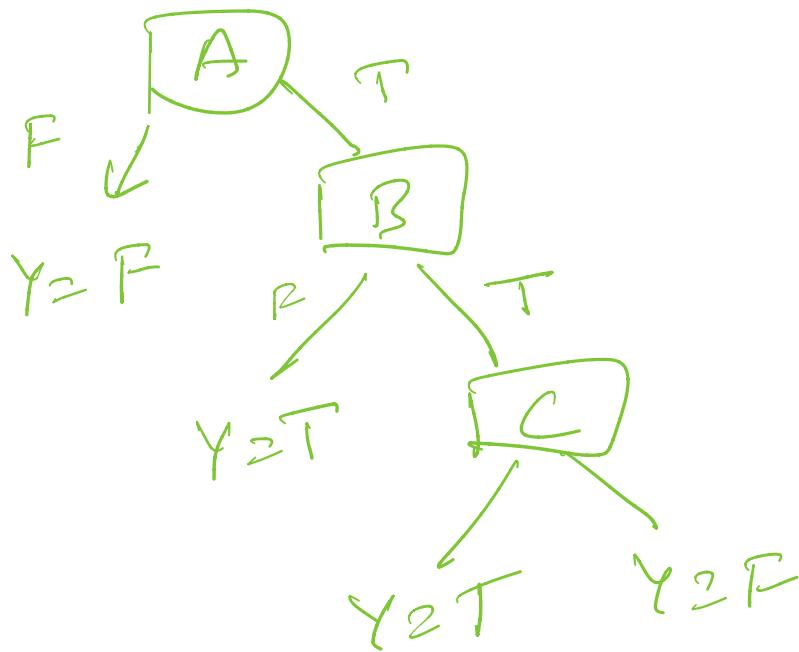
(Part 1) Draw the tree that would be learned by the greedy algorithm with zero training error. You do not need to show any computation.

A	B	C	Y
F	F	F	F
T	F	T	T
T	T	F	T
T	T	T	F

(Part 2) Is this tree optimal (i.e. does it get zero training error with minimal depth)? Explain in less

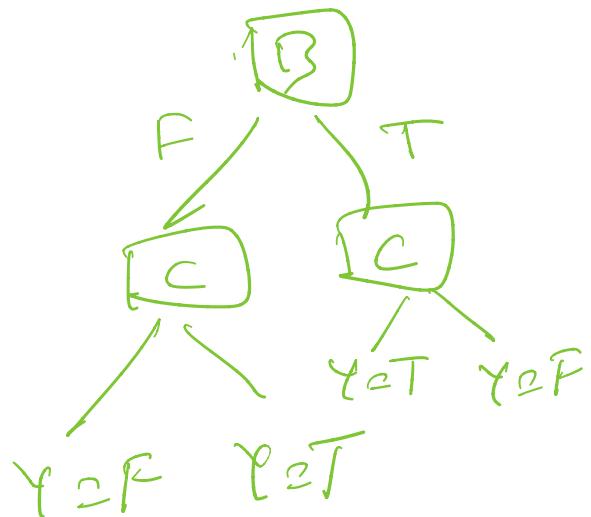
than two sentences. If it is not optimal, draw the optimal tree as well.

## Part 1



## Part II

Note:  $Y = B \otimes_C C$   
 $\therefore$  Depth 2 tree



(d) (4 points) Compute the dual of the following problem. Minimize  $x^2 + 1$  such that  $(x - 4)(x - 12) \leq 10$ .

$$\text{Ans: } x^2 - 16x + 38 \leq 0$$

$$\begin{aligned} L(x, \gamma) &= x^2 + 1 + \gamma(x^2 - 16x + 38) \\ &= x^2(\gamma + 1) - 16\gamma x + 38\gamma + 1 \end{aligned}$$

$$\text{Dual: } G(\gamma) = \min_x L(x, \gamma)$$

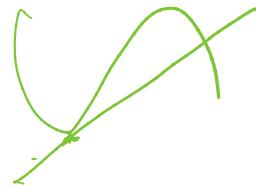
$$L(n, \lambda) = n^2(\lambda+1) - 16\lambda n + 38\lambda + 1$$

$$\Rightarrow \frac{\partial L}{\partial n} = 2n(\lambda+1) - 16\lambda = 0$$

$$\Rightarrow n = \frac{8\lambda}{\lambda+1}$$

$$L(\lambda) = \frac{64\lambda^2}{\lambda+1} - \frac{128\lambda^2}{\lambda+1} + \frac{304\lambda^2}{\lambda+1} + 1$$

Drd: met  $L(\lambda)$   
 $\lambda \geq 0$



### Question 2: True or False with Explanations

[8 pts] Each question below is for 2 points each. Please do not just write true or false. You also need to provide explanations. Just true or false will yield no points.

- (a) Given a function  $f$ , the set of sub-gradients is always non-empty.

False. Set of subgrads is non-empty only for Convex

- (b)  $k$ -nearest neighbor and decision tree models can be used for regression.

True. Reg: Take avg of  $j$  preds

- (c) No classifier can do better than a naive Bayes classifier if the distribution of the data is known.

False: NB is restricted to of cond Ind logistic mln reg loss

- (d) Maximizing the likelihood of linear regression yields multiple local optimums

True: Single local optima  
only if LR is strictly convex  
 $\Rightarrow$  Data matrix is PSD.

Convex



Global optima = Local optima

[Multiple local/global optima]



Strictly Convex



unique local/global optima.

### Question 3: Naive Bayes

[10 pts] Given the following training data points, provide the output of naive bayes classifier for test data points  $z_1$  and  $z_2$ :

$$\begin{aligned}x_1 &= (0, 0, 0, 1, 0, 0, 1), y_1 = 1 \\x_2 &= (0, 0, 1, 1, 0, 0, 0), y_2 = 1 \\x_3 &= (1, 1, 0, 0, 0, 1, 0), y_3 = -1 \\x_4 &= (1, 0, 0, 0, 1, 1, 0), y_4 = -1 \\x_5 &= (1, 1, 1, 1, 1, 1, 1), y_5 = 1 \\x_6 &= (0, 0, 0, 0, 0, 0, 0), y_6 = 1 \\x_7 &= (1, 1, 1, 1, 1, 1, 1), y_7 = -1 \\z_1 &= (1, 0, 0, 0, 0, 1, 0) \\z_2 &= (0, 1, 1, 0, 0, 1, 1)\end{aligned}$$

$$P(y=1) = \frac{4}{7}$$

$$P(y=-1) = \frac{3}{7}$$

$$P(m=1|y)$$

	na	nb	nc	nd	ne	nf	ng
y=1	1/4	1/4	2/4	2/4	1/4	1/4	1/4
y=-1	1	2/3	1/3	2/3	1/3	1	1/3

$$P(y=1|z_1) = \left[ \frac{1}{4} \cdot \frac{3}{4} \cdot \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{3}{4} \cdot \frac{1}{4} \cdot \frac{2}{4} \right]^{\frac{1}{7}}$$

$$\frac{36}{2^{12}} \cdot 7$$

$$P(y=-1|z_1) = \left[ \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{3}{3} \right]^{\frac{1}{7}}$$

$$= 24 / 3^7 \times 7$$

$$\frac{36}{2^{12} \times 7} < \frac{24}{3^5 \times 7}$$

$$8748 < 96304$$

$$z_1 : y = -1$$

For  $z_2$ , observe  $p(z_2 | y = -1) = 0$

$\therefore z_2$  classifies as  $y = 1$

( $p(z_2 | y = 1) > 0$ )

$$\omega_0 = \text{zero vector} \quad \lambda = 1$$

#### Question 4: Perceptrons

[10 pts] Demonstrate how the perceptron without bias (i.e. we set the parameter  $b = 0$  and keep it fixed) updates its parameters given the following training sequence:

$$\begin{aligned}x_1 &= (0, 0, 0, 1, 0, 0, 1), y_1 = 1 \\x_2 &= (1, 1, 0, 0, 0, 1, 0), y_2 = -1 \\x_3 &= (0, 0, 1, 1, 0, 0, 0), y_3 = 1 \\x_4 &= (1, 0, 0, 0, 1, 1, 0), y_4 = -1 \\x_5 &= (1, 0, 0, 0, 0, 1, 0), y_5 = -1\end{aligned}$$

Start  $w = (0, 0, 0, 0, 0, 0)$

$$y_1(w^T n_1) = 0 \leq 0 \rightarrow \text{yes}$$

$$w = w + y_1 n_1 = (0, 0, 0, 1, 0, 0, 1)$$

$$y_2(w^T n_2) = 0 \leq 0 \rightarrow \text{yes}$$

$$w = w + y_2 n_2$$

$$= (-1, -1, 0, 1, 0, -1, 1)$$

$$y_3(w^T n_3) = 1 \leq 0 \rightarrow \text{NO}$$

No update.

$$y_4 w^T x_4 = 2 \leq 0 \quad \text{No}$$

No update.

$$y_5 w^T x_5 = 2 \leq 0 \quad \text{No}$$

No update.

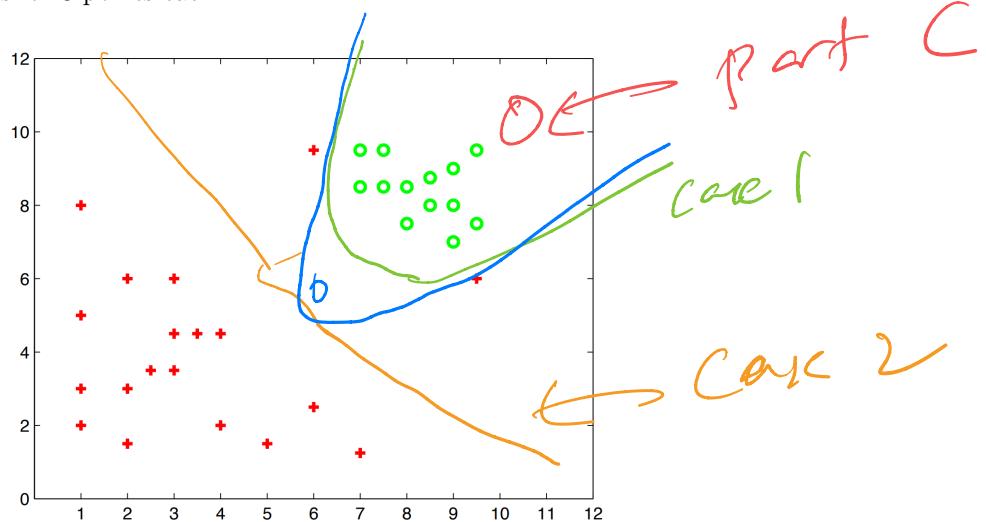
$$\underbrace{w = (-1, -1, 0, 1, 0, -1, 1)}$$

$$y_i(w^T x_i) \geq 0 \quad \forall i \Rightarrow \text{Train error} = 0$$

### Question 5: SVMs

[24 pts] This question will cover support vector machines.

- (a) (15 points) Given the following dataset (shown in the figure below), assume we are training the SVM with a quadratic kernel. The slack penalty  $C$  will determine the location of the separating hyper-plane. Each question below is for 3 points each.



- Where would the decision boundary be if  $C \rightarrow \infty$ ?

Perfectly separate Data

- For  $C \approx 0$ , where would the decision boundary be?

$C \approx 0$ : no penalty on misclass<sup>→</sup>

⇒ Maximize margin

- If we know that we cannot fully trust the obtained data points, which of the scenarios would we

prefer to train the model? (Large  $C$  or Small  $C$ )

Small  $C$  since let can be  
wrong (not reg<sup>n</sup>)

- Given  $C$  is large, draw an additional data point which *will not* change the decision boundary. Justify.

Drawn 0  
(not support vec)

- Given  $C$  is large, draw an additional data point which *will* change the decision boundary. Justify.

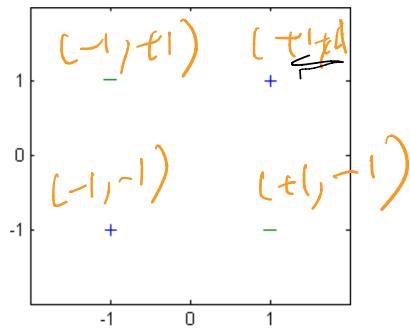
Drawn 0  
support vec : will change  
Dec<sup>n</sup> boundary

- (b) (9 points) Consider the data set below. Under which of the following feature vectors is the data linearly separable? For full credit, you must justify your answer by either providing a linear separator or explaining why such a separator does not exist. Each part is 1.8 points each

- $\phi(x_1, x_2) = [x_1 + x_2, x_1 - x_2]$
- $\phi(x_1, x_2) = [x_1^2, x_2^2, x_1 x_2]$
- $\phi(x_1, x_2) = [\exp(x_1), \exp(x_2)]$
- $\phi(x_1, x_2) = [x_1 \sin x_2, x_1]$
- $\phi(x_1, x_2) = [x_1 x_2, x_1]$

$\phi(n_1, n_2)$

$$\pi_1, \pi_2 = 0$$



①  $(\pi_1 + \pi_2, \pi_1 - \pi_2)$   $\nrightarrow$   
 will not separate  
 bcs orig. data is not LinSep

②  $(\pi_1^2, \pi_2^2, \pi_1 \pi_2)$   $\nrightarrow$

$$[0, 0, 0] (a_1^2 \pi_1^2 \pi_1 \pi_2) = 0$$

at dec. boundary

③  $(\exp(\pi_1), \exp(\pi_2)) \rightarrow$  ND

Since  $\pi_1 > y$  min fct.

$$(\underbrace{(\frac{1}{e}, \frac{1}{e}), (\frac{1}{e}, e), (e, \frac{1}{e})}_{10}, (e, e)) \in$$

④  $[n_1 \sin n_2, n_1] \rightarrow 4 \rightarrow$

since  $[n_1 \sin n_2, n_1] (k, 0) = 0$   
is sep hyperplane.

⑤  $[n_1 n_2, n_1] \rightarrow 4 \rightarrow$

since  $(n_1 n_2, n_1) [1, 0] = 0$   
is sep hyperplane

### Question 6: Logistic Regression

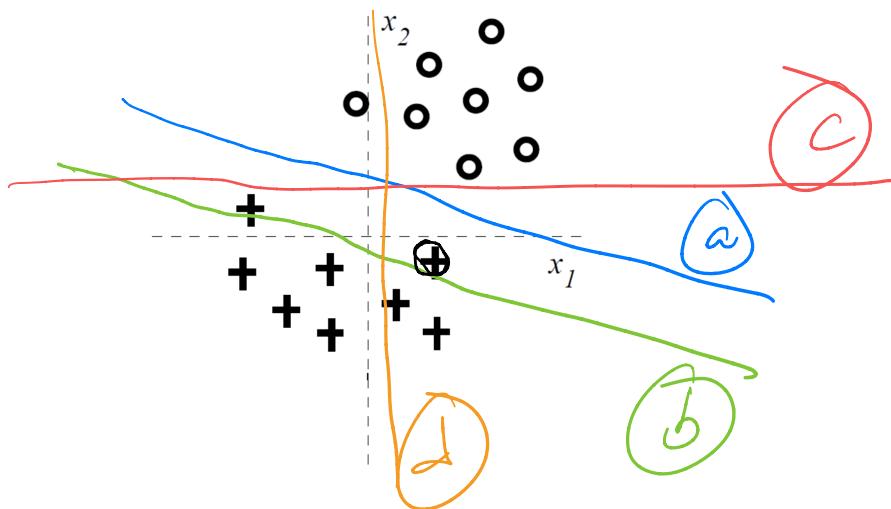
[18 pts] Given a two dimensional dataset shown in the figure below, we attempt to solve a simple binary classification using a linear logistic regression. The model is:

$$p(y=1|x, w_0, w_1, w_2) = \frac{1}{1 + \exp(-w_0 - w_1x_1 - w_2x_2)} \quad (2)$$

Consider training a regularized linear logistic regression model where we try to maximize:

$$\sum_{i=1}^n \log P(y_i|x_i, w_0, w_1, w_2) - C_0 w_0^2 - C_1 w_1^2 - C_2 w_2^2 \quad (3)$$

Note we have different regularization parameters for each coordinate. Draw the approximate decision boundary.



arises in the following cases (each sub-question is 3 points). Also explain what will be the training error in each case.

- (a) When  $C_0, C_1, C_2 \approx 0$

No reg : Train error  $\Rightarrow$

- (b) When  $C_0$  is very large and  $C_1, C_2 \approx 0$

No bias term : Train error  $\Rightarrow$   
 $\underbrace{\text{f. } \{ \}_{\text{--}} \text{--}}$

$$T, \varepsilon = 0$$

(c) When  $C_1$  is very large and  $C_0, C_2 \approx 0$

$C_1$  is large  $\Rightarrow w_1 \rightarrow 0$

$\therefore$  DB will depend less on  $w_1$  &  
non-convex

(d) When  $C_2$  is very large and  $C_0, C_1 \approx 0$

$C_2$  is large  $\Rightarrow w_2 \rightarrow 0$

DB will be vertical

$$T, \varepsilon > 0$$

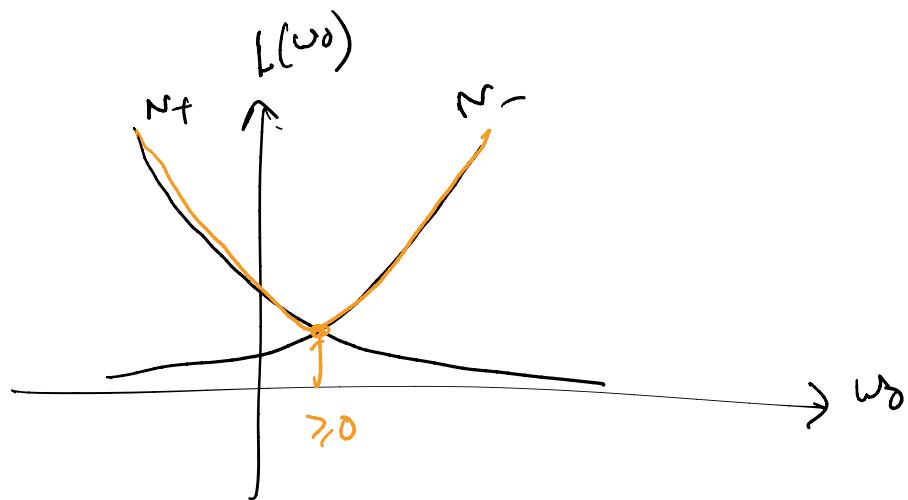
(e) Consider the case when  $C_1$  and  $C_2$  are very large and  $C_0 \approx 0$ . What is the value of  $w_0$  that we expect to obtain?

$$\min_{w_0} N \left[ \log((1 + \exp(w_0)) + \log(1 + \exp(-w_0))) \right]$$

(f) In the above case, assume we add a few more positives to this dataset (i.e. make it imbalanced). What is the value of  $w_0$  we expect? (you can give a range of values of  $w_0$  if you prefer).

$$\frac{N \log(1 + \exp(-w_0))}{\uparrow + m \log(1 + \exp(w_0))} \leftarrow m \uparrow$$

Dominating



$N_+$  Dominates over  $N_-$

so  $s_0^r$  will be  $\geq 0$

### Question 7: Decision Trees

[10 pts] Consider the dataset shown below. We will use this dataset to learn a decision tree which predicts if people pass machine learning (Yes or No), based on their previous GPA (High, Medium, or Low) and whether or not they studied. For this problem, assume that  $\log_2 3 \approx 1.6$  (it is ok if you leave the answers in terms of the logs as well).

GPA	Studied	Passed
L	F	F
L	T	T
M	F	F
M	T	T
H	F	T
H	T	T

terms of the logs as well).

- (a) (6 points) Compute  $H(\text{Passed})$ ,  $H(\text{Passed} | \text{GPA})$  and  $H(\text{Passed} | \text{Studied})$

$$\begin{aligned}
 H(\text{Passed}) &= -\frac{2}{6} \log_2 \frac{2}{6} - \frac{4}{6} \log_2 \frac{4}{6} \\
 &= -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{2} \\
 &\approx \log_2 3 - \frac{2}{3} \approx 0.92
 \end{aligned}$$

$$\text{Sim. } H(\text{Passed} | \text{GPA}) = \frac{2}{3} \approx 0.66$$

$$\begin{aligned}
 H(\text{Passed} | \text{Studied}) &= \frac{1}{2} \log_2 3 - \frac{1}{3} \\
 &\approx 0.46
 \end{aligned}$$

(b) (4 points) Draw the full decision tree that would be learned for this dataset.

